



**HAL**  
open science

## A kernel-based approach to non-stationary reinforcement learning in metric spaces

Omar D Domingues, Pierre Ménard, Matteo Pirotta, Emilie Kaufmann,  
Michal Valko

► **To cite this version:**

Omar D Domingues, Pierre Ménard, Matteo Pirotta, Emilie Kaufmann, Michal Valko. A kernel-based approach to non-stationary reinforcement learning in metric spaces. International Conference on Artificial Intelligence and Statistics, Apr 2021, San Diego / Virtual, United States. hal-03289026

**HAL Id: hal-03289026**

**<https://inria.hal.science/hal-03289026v1>**

Submitted on 16 Jul 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# A Kernel-Based Approach to Non-Stationary Reinforcement Learning in Metric Spaces

---

Omar D. Domingues<sup>1,2</sup> Pierre Ménard<sup>3</sup> Matteo Pirotta<sup>4</sup> Emilie Kaufmann<sup>1,2,5</sup> Michal Valko<sup>1,6</sup>  
<sup>1</sup>Inria Lille <sup>2</sup>Université de Lille <sup>3</sup>OvGU <sup>4</sup>Facebook AI Research <sup>5</sup>CNRS <sup>6</sup>DeepMind Paris

## Abstract

In this work, we propose **KeRNS**: an algorithm for episodic reinforcement learning in non-stationary Markov Decision Processes (MDPs) whose state-action set is endowed with a metric. Using a non-parametric model of the MDP built with time-dependent kernels, we prove a regret bound that scales with the covering dimension of the state-action space and the total variation of the MDP with time, which quantifies its level of non-stationarity. Our method generalizes previous approaches based on sliding windows and exponential discounting used to handle changing environments. We further propose a practical implementation of **KeRNS**, we analyze its regret and validate it experimentally.

## 1 Introduction

In reinforcement learning (RL), an agent interacts with an environment by sequentially taking actions, receiving rewards and observing state transitions. One of the main challenges in RL is the trade-off between exploration, the act of gathering information about the environment, and exploitation, the act of using the current knowledge to maximize the sum of rewards. In non-stationary environments, handling this trade-off becomes much harder: what has been learned in the past may no longer be valid in the present. Therefore, the agent needs to constantly re-explore previously known parts of the environment to discover possible changes. In this work, we propose **KeRNS**,<sup>1</sup> an algorithm that handles this problem by acting optimistically and by

---

<sup>1</sup>meaning Kernel-based Reinforcement Learning in Non-Stationary environments.

forgetting data that are far in the past, which naturally causes the agent to keep exploring to discover changes. **KeRNS** relies on non-parametric kernel estimators of the MDP, and the non-stationarity is handled by using time-dependent kernels.

The regret of an algorithm, defined as the difference between the rewards obtained by an optimal agent and the ones obtained by the algorithm, allows us to quantify how well an agent balances exploration and exploitation. We prove a regret bound for **KeRNS** that holds in a challenging setting, where the state-action space can be continuous and the environment can change in every episode, as long as the cumulative changes remain small when compared to the total number of episodes.

**Related work** Regret bounds for RL in stationary environments have been extensively studied in finite (tabular) MDPs (Jaksch et al., 2010; Azar et al., 2017; Dann et al., 2017; Jin et al., 2018; Zanette and Brunskill, 2019), and also in metric spaces under Lipschitz continuity assumptions (Ortner and Ryabko, 2012; Song and Sun, 2019; Sinclair et al., 2019; Domingues et al., 2020; Sinclair et al., 2020). Recent works provide algorithms with regret bounds for non-stationary RL in the tabular setting (Gajane et al., 2018; Ortner et al., 2019; Cheung et al., 2020). These algorithms estimate the transitions and the rewards in an episode  $k$  using the data observed up to episode  $k - 1$ . However, since the MDP can change from one episode to another, these estimators are *biased*. If nothing is done to handle this bias, the algorithms will suffer a linear regret (Ortner et al., 2019) that depends on the magnitude of the bias. To deal with this issue, different approaches have been proposed: Gajane et al. (2018) and Cheung et al. (2020) use sliding windows to compute estimators that use only the most recently observed transitions, whereas Ortner et al. (2019) restart the algorithm periodically and, after each restart, new estimators are build and past data are discarded. In the multi-armed bandit literature, in addition to sliding windows, exponential discounting has also been used as a mean to give more importance to recent data (Kocsis and Szepesvári, 2006;

Garivier and Moulines, 2011; Russac et al., 2019). In this paper, we study the *dynamic regret* of the algorithm, where, in each episode  $k$ , we compare the learner to the optimal policy of the MDP in episode  $k$ . A related approach consists in comparing the performance of the learner to the best stationary policy in hindsight, e.g., (Even-Dar et al., 2009; Yu and Mannor, 2009; Neu et al., 2013; Dick et al., 2014), which is less suited to non-stationary environments, since the performance of any fixed policy can be very bad. Non-stationary RL has also been studied outside the regret minimization framework, without, however, tackling the issue of exploration. For instance, Choi et al. (2000) propose a model where the MDP varies according to a sequence of tasks whose changes form a Markov chain. Szita et al. (2002) and Csáji and Monostori (2008) study the convergence of Q-learning when the environment changes but remain close to a fixed MDP. Assuming full knowledge of the MDP at each time step, but with unknown evolution, Lecarpentier and Rachelson (2019) introduce a risk-averse approach to planning in slowly changing environments. In a related setting, Lykouris et al. (2019) study episodic RL problems where the MDP can be corrupted by an adversary and provide regret bounds in this case.

**Contributions** We provide the first regret bound for non-stationary RL in continuous environments. More precisely, we show that the **Kernel-UCBVI** algorithm of Domingues et al. (2020), based on non-parametric kernel smoothing, can be modified to tackle non-stationary environments by using appropriate time- and space-dependent kernels. We analyze the resulting algorithm, **KeRNS**, under mild assumptions on the kernel, which in particular recover previously studied forgetting mechanisms to tackle non-stationarity in bandits and RL: sliding windows (Gajane et al., 2018) and exponential discounting (Kocsis and Szepesvári, 2006; Garivier and Moulines, 2011; Russac et al., 2019), and allow for combinations between those. On the practical side, kernel-based approaches can be very computationally demanding since their complexity grows with the number of data points. Building on the notion of representative states, promoted in previous work on practical kernel-based RL (Kveton and Theodorou, 2012; Barreto et al., 2016) we propose an efficient version of **KeRNS**, called **RS-KeRNS**, which has constant runtime per episode. We analyze the regret of **RS-KeRNS**, showing that it enables a trade-off between regret and runtime, and we validate this algorithm empirically.

## 2 Setting

**Notation** For any  $n \in \mathbb{N}^*$ , let  $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ . If  $\mu$  and  $P(\cdot|x, a)$  are measures for any  $(x, a)$  and  $f$  is an

arbitrary function, we define  $\mu f \stackrel{\text{def}}{=} \int f(y) d\mu(y)$  and  $Pf(x, a) \stackrel{\text{def}}{=} \int f(y) dP(y|x, a)$ .<sup>2</sup>

**Non-stationary MDPs** We consider an episodic RL setting where, in each episode  $k \in [K]$ , an agent interacts with the environment for  $H \in \mathbb{N}^*$  time steps. The time is indexed by  $(k, h)$ , where  $k$  represents an episode and  $h$  the time step within the episode. The environment is modeled as a non-stationary MDP, defined by the tuple  $(\mathcal{X}, \mathcal{A}, r, P)$ , where  $\mathcal{X}$  is the state space,  $\mathcal{A}$  is the action space,  $r = \{r_h^k\}_{k,h}$  and  $P = \{P_h^k\}_{k,h}$  are sets of reward functions and transition kernels, respectively. More precisely, when taking action  $a$  in state  $x$  at time  $(k, h)$ , the agent observes a random reward  $\tilde{r}_h^k \in [0, 1]$  with mean  $r_h^k(x, a)$  and makes a transition to the next state according to the probability measure  $P_h^k(\cdot|x, a)$ . A deterministic policy  $\pi$  is a mapping from  $[H] \times \mathcal{X}$  to  $\mathcal{A}$ , and we denote by  $\pi(h, x)$  the action chosen in state  $x$  at step  $h$ . The action-value function of a policy  $\pi$  in step  $h$  of episode  $k$  is defined as

$$Q_{k,h}^\pi(x, a) \stackrel{\text{def}}{=} \mathbb{E} \left[ \sum_{h'=h}^H r_{h'}^k(x_{h'}, a_{h'}) \middle| x_h = x, a_h = a \right]$$

where  $x_{h'+1} \sim P_{h'}^k(\cdot|x_{h'}, a_{h'})$ ,  $a_{h'} = \pi(h', x)$ , and its value function is defined by  $V_{k,h}^\pi(x) = Q_{k,h}^\pi(x, \pi(h, x))$ . The optimal value functions,  $V_{k,h}^* \stackrel{\text{def}}{=} \sup_\pi V_{k,h}^\pi(x)$  satisfy the Bellman equations (Puterman, 2014)

$$V_{k,h}^*(x) = \max_{a \in \mathcal{A}} Q_{k,h}^*(x, a), \text{ where} \\ Q_{k,h}^*(x, a) \stackrel{\text{def}}{=} r_h^k(x, a) + P_h^k V_{k,h+1}^*(x, a)$$

and where  $V_{k,H+1}^* = 0$  by definition.

**Dynamic regret** The agent interacts with the environment in a sequence of episodes and, in each episode  $k$ , it uses a policy  $\pi_k$  that can be chosen based on its observations from previous episodes. We measure its performance by the dynamic regret, defined as the sum over all episodes of the difference between the optimal value function in episode  $k$  and the value of  $\pi_k$ :

$$\mathcal{R}(K) \stackrel{\text{def}}{=} \sum_{k=1}^K \left( V_{k,1}^*(x_1^k) - V_{k,1}^{\pi_k}(x_1^k) \right)$$

where  $x_1^k$  is the starting state in each episode, which is chosen arbitrarily and given to the learner.

**Assumptions** Since regret lower bounds scale with the number of states and actions (Jaksch et al., 2010), structural assumptions are needed in order to enable

<sup>2</sup>See also Table 2 in Appendix A summarizing the main notations used in the paper and in the proofs.

learning in continuous MDPs. A common assumption is that rewards and transitions are Lipschitz continuous with respect to some known metric (Ortner and Ryabko, 2012; Song and Sun, 2019; Domingues et al., 2020; Sinclair et al., 2020), which is the approach that we follow in this work. We make no assumptions regarding how the MDP changes, and our regret bounds will be expressed in terms of its total variation over time.

**Assumption 1.** *The state-action space  $\mathcal{X} \times \mathcal{A}$  is equipped with a metric  $\rho : (\mathcal{X} \times \mathcal{A})^2 \rightarrow \mathbb{R}_+$ , which is given to the learner. Also, we assume that there exists a metric  $\rho_{\mathcal{X}}$  on  $\mathcal{X}$  such that, for all  $(x, x', a)$ ,  $\rho[(x, a), (x', a)] \leq \rho_{\mathcal{X}}(x, x')$ .<sup>3</sup>*

**Assumption 2.** *The reward functions are  $L_r$ -Lipschitz and the transition kernels are  $L_p$ -Lipschitz with respect to the 1-Wasserstein distance:  $\forall(x, a, x', a') \text{ and } \forall(k, h) \in [K] \times [H]$ ,*

$$|r_h^k(x, a) - r_h^k(x', a')| \leq L_r \rho[(x, a), (x', a')], \text{ and} \\ \mathbb{W}_1(P_h^k(\cdot|x, a), P_h^k(\cdot|x', a')) \leq L_p \rho[(x, a), (x', a')]$$

where, for two measures  $\mu$  and  $\nu$ , we have  $\mathbb{W}_1(\mu, \nu) \stackrel{\text{def}}{=} \sup_{f: \text{Lip}(f) \leq 1} \int_{\mathcal{X}} f(y) (d\mu(y) - d\nu(y))$  and where, for any Lipschitz function  $f : \mathcal{X} \rightarrow \mathbb{R}$  with respect to  $\rho_{\mathcal{X}}$ ,  $\text{Lip}(f)$  denotes its Lipschitz constant.

**Assumption 3.** *For any  $(k, h)$ , the optimal  $Q$ -function  $Q_{k,h}^*$  is  $L$ -Lipschitz with respect to  $\rho$ . Assumptions 1 and 2 imply that  $L \leq \sum_{h=1}^H L_r L_p^{H-h}$  (Lemma 26 in the Appendix).*

### 3 An Algorithm for Kernel-Based RL in Non-Stationary Environments

In this section, we introduce **KeRNS**, a model-based RL algorithm for learning in non-stationary MDPs. In each episode  $k$ , we estimate the transitions and the rewards using the data observed up to episode  $k - 1$ . Using exploration bonuses that represent the uncertainty in the estimated model, **KeRNS** builds a  $Q$ -function  $Q_h^k$ , and plays the greedy policy with respect to it. **KeRNS** generalizes sliding-window and exponential discounting approaches by considering time-dependent kernel functions, which also allow us to handle exploration in continuous environments (Domingues et al., 2020).

#### 3.1 Kernel-Based Estimators for Changing MDPs

Let  $\Gamma : \mathbb{N} \times (\mathcal{X} \times \mathcal{A})^2 \rightarrow [0, 1]$  be a *non-stationary kernel function*, where  $\Gamma(t, u, v)$  represents the similarity

<sup>3</sup>If  $(\mathcal{A}, \rho_{\mathcal{A}})$  is also a metric space, we can take  $\rho[(x, a), (x', a')] = \rho_{\mathcal{X}}(x, x') + \rho_{\mathcal{A}}(a, a')$ , for instance. See Section 2.3 of Sinclair et al. (2019) for more examples and a discussion.

between two state action pairs  $u, v$  in  $\mathcal{X} \times \mathcal{A}$  visited at an interval  $t$ .

**Definition 1** (kernel weights). *Let  $(x_h^s, a_h^s)$  be the state-action pair visited at time  $(s, h)$ . For any  $(x, a) \in \mathcal{X} \times \mathcal{A}$  and  $s < k$ , we define the weights and the normalized weights at time  $(k, h)$  as*

$$w_h^{k,s}(x, a) \stackrel{\text{def}}{=} \Gamma(k - s - 1, (x, a), (x_h^s, a_h^s))$$

and  $\tilde{w}_h^{k,s}(x, a) \stackrel{\text{def}}{=} w_h^{k,s}(x, a) / \mathbf{C}_h^k(x, a)$ , where  $\mathbf{C}_h^k(x, a) \stackrel{\text{def}}{=} \beta + \sum_{s=1}^{k-1} w_h^{k,s}(x, a)$  and  $\beta > 0$  is a regularization parameter.

Using the kernel function  $\Gamma$  and past data, **KeRNS** builds estimators  $\hat{r}_h^k$  of the reward function and  $\hat{P}_h^k$  of the transitions at time  $(k, h)$ , which are defined below.

**Definition 2** (empirical MDP). *At time  $(s, h) \in [K] \times [H]$ , let  $(x_h^s, a_h^s, x_{h+1}^s, \tilde{r}_h^s)$  represent the state, the action, the next state and the reward observed by the algorithm. Before each episode  $k$ , **KeRNS** estimates the rewards and transitions using the data observed up to episode  $k - 1$ :*

$$\hat{r}_h^k(x, a) \stackrel{\text{def}}{=} \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) \tilde{r}_h^s, \\ \hat{P}_h^k(y|x, a) \stackrel{\text{def}}{=} \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) \delta_{x_{h+1}^s}(y)$$

where  $\delta_x$  is the Dirac measure at  $x$ . Let  $\widehat{\mathcal{M}}_k$  be the MDP whose rewards and transitions at step  $h$  are  $\hat{r}_h^k(x, a)$  and  $\hat{P}_h^k(y|x, a)$ .<sup>4</sup>

The weights  $w_h^{k,s}(x, a)$  measure the influence that the transitions and rewards observed at time  $(s, h)$  will have on the estimators for the state-action pair  $(x, a)$  at time  $(k, h)$ . Their sum,  $\mathbf{C}_h^k(x, a)$ , is a proxy for the number of visits to  $(x, a)$ . Intuitively, the kernel function  $\Gamma$  must be designed in order to ensure that  $w_h^{k,s}(x, a)$  is small when  $(x, a)$  is very far from  $(x_h^s, a_h^s)$ , with respect to the distance  $\rho$ . It must also be small when  $k - s - 1$  is large, which means that the sample  $(x_h^s, a_h^s)$  was collected too far in the past and should have a small impact on the estimators. For our theoretical analysis, we will need the assumptions below on the kernel function  $\Gamma$ .

**Assumption 4** (kernel properties). *Let  $\sigma > 0$ ,  $\eta \in ]0, 1[$  and  $W \in \mathbb{N}$  be the kernel parameters. For each set of parameters, we assume that we have access to a base kernel function  $\bar{\Gamma}_{(\eta, W)} : \mathbb{N} \times \mathbb{R} \rightarrow [0, 1]$  and we define, for any  $t, u, v \in \mathbb{N}^* \times \mathcal{X} \times \mathcal{A}$ ,*

$$\Gamma(t, u, v) = \bar{\Gamma}_{(\eta, W)}(t, \rho[u, v] / \sigma).$$

<sup>4</sup>Since the normalized weights do not sum to 1,  $\hat{P}_h^k$  is not a probability kernel. In this case, we suffer a bias of order  $\beta$  and the property that  $\hat{P}_h^k$  is a sub-probability measure is enough for the analysis.

We assume that  $z \mapsto \bar{\Gamma}_{(\eta, W)}(t, z)$  is non-increasing for any  $t \in \mathbb{N}$ . Additionally, we assume that there exists positive constants  $C_1, C_2$ , a constant  $C_3 \geq 0$  and an arbitrary function  $G : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  that satisfies  $G(4) > 0$  such that

- (1)  $\forall(t, z), \bar{\Gamma}_{(\eta, W)}(t, z) \leq C_1 \exp(-z^2/2)$
- (2)  $\forall(t, y, z), |\bar{\Gamma}_{(\eta, W)}(t, y) - \bar{\Gamma}_{(\eta, W)}(t, z)| \leq C_2 |y - z|$
- (3)  $\forall z, \bar{\Gamma}_{(\eta, W)}(t, z) \leq C_3 \eta^t$ , for all  $t \geq W$
- (4)  $\forall z, \bar{\Gamma}_{(\eta, W)}(t, z) \geq G(z) \eta^t$ , for all  $t < W$ .

We now provide some justification for these conditions. (1) ensures that the bias due to kernel smoothing remains bounded by  $\tilde{\mathcal{O}}(\sigma)$  (Lemma 23); (2) ensures smoothness conditions that are needed to provide concentration inequalities for the rewards and transitions (Lemma 24); (3) and (4) allow us to control the bias and the variance due to non-stationarity, respectively (Lemmas 2 and 16). Intuitively, (3) says the algorithm should forget data further than  $W$  episodes in the past, and (4) says that recent data in the  $W$  most recent episodes must have a minimum weight. The condition  $G(4) > 0$  is mostly technical: it is used to ensure that  $\mathbf{C}_h^k(x, a)$  is not too small in a  $4\sigma$ -neighborhood of  $(x, a)$  (see lemmas 15 and 16). The kernels in the example below satisfy our conditions, and show that they indeed generalize sliding-window and exponential discounting approaches:

**Example 1** (sliding-window and exponential discount). The kernels  $\bar{\Gamma}_{(\eta, W)}(t, z) = \mathbb{I}\{t < W\} \exp(-|z|^p/2)$  (sliding-window) and  $\bar{\Gamma}_{(\eta, W)}(t, z) = \eta^t \exp(-|z|^p/2)$  (exponential discount) satisfy Assumption 4 for  $p \geq 2$ .

The conditions in Assumption 4 are needed to prove our regret bounds. However, if one has further knowledge about the MDP and its changes, this information can also be integrated to the kernel function  $\Gamma$ . For example, if the MDP only changes in certain region of the state-action space, the kernel can be designed to forget past data only in that region. Also, the kernel  $\Gamma$  can be designed to enforce restarts, as proposed by Ortner et al. (2019) for finite MDPs, by setting  $\Gamma(t, u, v)$  to zero every time  $t$  exceeds a certain threshold. Although this would require a separate analysis, our proof could be combined to the one of (Ortner et al., 2019) to obtain a regret bound in this case.

### 3.2 Algorithm

**KeRNS** is presented in Algorithm 1. At time  $(k, h)$ , let  $\mathbf{B}_h^k(x, a)$  be the exploration bonus at  $(x, a)$  representing

the uncertainty of  $\widehat{\mathcal{M}}_k$  with respect to the true MDP:

$$\mathbf{B}_h^k(x, a) = \tilde{\mathcal{O}} \left( \frac{H}{\sqrt{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} + L\sigma \right) \quad (1)$$

where  $\tilde{\mathcal{O}}(\cdot)$  hides logarithmic terms. The exact expression for the bonuses is given in Def. 5 in Appendix A. Before starting episode  $k$ , **KeRNS** computes, for all  $h \in [H]$ , the values  $Q_h^k$  by running backward induction on  $\widehat{\mathcal{M}}_k$ , with the bonus  $\mathbf{B}_h^k(x, a)$  added to the rewards, followed by an interpolation step:

$$\begin{aligned} \tilde{Q}_h^k(x, a) &= \hat{r}_h^k(x, a) + \hat{P}_h^k V_{h+1}^k(x, a) + \mathbf{B}_h^k(x, a) \\ Q_h^k(x, a) &= \min_{s \in [k-1]} \left( \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho[(x, a), (x_h^s, a_h^s)] \right) \\ V_h^k(x) &= \min(H - h + 1, \max_a Q_h^k(x, a)) \end{aligned}$$

where  $V_{H+1}^k \stackrel{\text{def}}{=} 0$ . The interpolation is needed to ensure that  $Q_h^k$  and  $V_h^k$  are  $L$ -Lipschitz. This procedure is defined in detail in Algorithm 3 in Appendix A, which is the same kind of backward induction used by **Kernel-UCBVI** (Domingues et al., 2020). Once  $Q_h^k$  is computed, **KeRNS** plays the greedy policy associated to it. Notice that, although  $Q_h^k(x, a)$  and  $V_h^k(x)$  are defined for all  $(x, a)$ , they only need to be computed for the states and actions observed by the algorithm up to episode  $k$ .

---

#### Algorithm 1 KeRNS

---

- 1: **Input:**  $K, H, L, L_r, L_p, \beta, \delta, d, \sigma, \eta, W$ .
  - 2: Initialize history:  $\mathcal{T}_h = \emptyset$  for all  $h \in [H]$ .
  - 3: **for** episode  $k = 1, \dots, K$  **do**
  - 4:   get initial state  $x_1^k$
  - 5:   // Run kernel backward induction
  - 6:   compute  $(Q_h^k)_h$  using  $(\mathcal{T}_h)_h$  and Algorithm 3.
  - 7:   **for**  $h = 1, \dots, H$  **do**
  - 8:     execute  $a_h^k = \arg\max_a Q_h^k(x_h^k, a)$
  - 9:     observe reward  $\hat{r}_h^k$  and next state  $x_{h+1}^k$
  - 10:     store transition  $\mathcal{T}_h = \mathcal{T}_h \cup \{x_h^k, a_h^k, x_{h+1}^k, \hat{r}_h^k\}$
  - 11:   **end for**
  - 12: **end for**
- 

### 3.3 Theoretical guarantees

We introduce  $\Delta$ , the total variation of the MDP in  $K$  episodes:

**Definition 3** (MDP variation). We define  $\Delta = \Delta^r + L\Delta^p$ , where

$$\begin{aligned} \Delta^r &\stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x, a} |r_h^i(x, a) - r_h^{i+1}(x, a)|, \\ \Delta^p &\stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x, a} \mathbb{W}_1(P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)) \end{aligned}$$

A similar notion has been introduced, for instance, by Ortner et al. (2019); Li and Li (2019) for MDPs and by Besbes et al. (2014) for multi-armed bandits. Here, the difference is that we use the Wasserstein distance to define the variation of the transitions, instead of the total variation (TV) distance  $\|P_h^i(\cdot|x, a) - P_h^{i+1}(\cdot|x, a)\|_1$ . This choice was made in order to take into account the metric  $\rho$  when measuring changes in the environment: our results would be analogous if we had chosen the TV distance.<sup>5</sup>

Using the same algorithm, we provide two regret bounds for **KeRNS**, which are given below. The notation  $\lesssim$  omits constants and logarithmic terms (see Definition 4 in Appendix A).

**Theorem 1.** *The regret of KeRNS is bounded as  $\mathcal{R}^{\text{KeRNS}}(K) \lesssim \min(\mathcal{R}_1(K), \mathcal{R}_2(K)) + \text{bias}(\sigma, \eta, W, \Delta)$ , where*

$$\begin{aligned} \mathcal{R}_1(K) &= H^2 K \sqrt{\log \frac{1}{\eta}} \sqrt{|\mathcal{C}'_\sigma| |\mathcal{C}_\sigma|} + H^2 |\mathcal{C}_\sigma| K \log \frac{1}{\eta} \\ \mathcal{R}_2(K) &= H^2 K \sqrt{\log \frac{1}{\eta}} \sqrt{|\mathcal{C}_\sigma|} + H^3 |\mathcal{C}_\sigma| |\mathcal{C}'_\sigma| K \log \frac{1}{\eta} \\ \text{bias}(\sigma, \eta, W, \Delta) &= W \Delta H + \frac{\eta^W}{1 - \eta} K H^3 + L K H \sigma \end{aligned}$$

with probability at least  $1 - \delta$ . Here,  $|\mathcal{C}'_\sigma|$  and  $|\mathcal{C}_\sigma|$  are the  $\sigma$ -covering numbers of  $(\mathcal{X}, \rho_{\mathcal{X}})$  and  $(\mathcal{X} \times \mathcal{A}, \rho)$  respectively,  $(\sigma, \eta, W)$  are the kernel parameters.

*Proof.* This result comes from combining theorems 3 and 4 in Appendix F. See Section 5 for a proof outline.

As discussed below, after optimizing the kernel parameters (Table 1), the bound  $\mathcal{R}_1$  has a worse dependence on  $K$ , and a better dependence on  $\Delta$ . On the other hand,  $\mathcal{R}_2$  is better with respect to  $K$ , but worse in  $\Delta$ . Concretely, this trade-off may give hints on how to choose the kernel parameters according to the amount of variation that we expect to see in the environment. Technically, the difference comes from how we handle the concentration of the transitions in the proof. To obtain  $\mathcal{R}_1$ , we use concentration inequalities on the term  $|(\hat{P}_h^k - P_h^k)f|$  for *all* functions  $f$  that are bounded and Lipschitz continuous. To obtain  $\mathcal{R}_2$ , the concentration is done only for  $f = V_{k, h+1}^*$ , but this results in larger second-order terms, as in (Azar et al., 2017; Domingues et al., 2020).

**Corollary 1.** *Let  $d$  be the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ . By optimizing the kernel parameters, we obtain the regret bounds in Table 1. Table 3 in Appendix B.2 gives the values of  $(\sigma, \eta, W)$  that yield these bounds.*

*Proof.* Assuming that  $|\mathcal{C}'_\sigma| \leq |\mathcal{C}_\sigma|$ , we have that  $|\mathcal{C}_\sigma|$  and  $|\mathcal{C}'_\sigma|$  are  $\mathcal{O}(1/\sigma^d)$ . Then, the bounds follow from

<sup>5</sup>More precisely, in the proof of Corollary 2, the Wasserstein distance could be replaced by the TV distance.

Theorem 1. The general case, handling separately the covering dimensions of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and  $(\mathcal{X}, \rho_{\mathcal{X}})$ , is stated in corollaries 6 and 9 in Appendix F.

**Discussion** We now discuss regret bounds for optimized kernel parameters, according to the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ , denoted by  $d$ . Roughly, the covering dimension is the smallest number  $d \geq 0$  such that the  $\sigma$ -covering number  $|\mathcal{C}_\sigma|$  is  $\mathcal{O}(1/\sigma^d)$ .<sup>6</sup> We consider two cases: the tabular (finite MDP) case, where the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$  is  $d = 0$ , and the continuous case, where  $d > 0$ .

**Tabular case** Let  $X = |\mathcal{X}|$  and  $A = |\mathcal{A}|$ . By taking  $\sigma = 0$ , we have  $|\mathcal{C}'_\sigma| = X$  and  $|\mathcal{C}_\sigma| = XA$ . As shown in Table 1, the  $\mathcal{R}_1$  bound states that the regret of **KeRNS** is  $\tilde{\mathcal{O}}\left(H^2 X \sqrt{A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}}\right)$ . This bound matches the one proved by Ortner et al. (2019) for the average reward setting using restarts, up to a factor of  $H^{\frac{2}{3}}$  coming from our episodic setting, where the transitions  $P_h^k$  depend on  $h$ . The  $\mathcal{R}_2$  bound states that the regret of **KeRNS** can be improved to  $\tilde{\mathcal{O}}\left(H^2 \sqrt{XA} \Delta^{\frac{1}{3}} K^{\frac{2}{3}}\right)$ , up to second-order terms. In the bandit case ( $H = 1$ ), these bounds are *optimal* in terms of  $K$  and  $\Delta$  (Besbes et al., 2014).

**Continuous case** For  $d > 0$ , we prove the first dynamic regret bounds in our setting, which are of order  $H^2 \Delta^{\frac{1}{3}} K^{\frac{2d+2}{2d+3}}$  (better in  $\Delta$ ) or  $H^2 \Delta^{\frac{1}{2}} K^{\frac{2d+1}{2d+2}}$  (better in  $K$ ) for two different tunings of the kernel. Deriving a lower bound in the non-stationary case for  $d > 0$  is an open problem, even for multi-armed bandits. As a sanity-check, we note that in stationary MDPs, for which  $\Delta = 0$ , we recover the regret bound of **Kernel-UCBVI**<sup>7</sup> (Domingues et al., 2020) of  $H^3 K^{\frac{2d}{2d+1}}$  from the bound  $\mathcal{R}_2$  with  $\log(1/\eta) = 1/K$ ,  $W \rightarrow \infty$  and  $\sigma = K^{-\frac{1}{2d+1}}$ , which is optimal for  $d = 1$  in the (stationary) bandit case (Bubeck et al., 2011).

In tabular MDPs, we may achieve sub-linear regret as long as  $\Delta < K$ .<sup>8</sup> In the continuous case however, our bounds show that we might need  $\Delta < K^{\frac{3}{2d+3}}$  (for the  $\mathcal{R}_1$  bound) or  $\Delta < K^{\frac{1}{d+1}}$  (for the  $\mathcal{R}_2$  bound) in order to avoid a linear regret, which is an immediate consequence of the bounds in Table 1.

**Knowledge of  $\Delta$**  To optimally choose the kernel parameters, **KeRNS** requires an upper bound on the

<sup>6</sup>For more details about covering numbers and covering dimension, see Section 3 of Kleinberg et al. (2019) and Section 2.2 of Sinclair et al. (2019).

<sup>7</sup>Another choice of  $\eta$  might allow us to avoid the dependence on  $H^3$  of **Kernel-UCBVI** and get  $H^2$  instead.

<sup>8</sup>Notice that, if  $\Delta$  scales linearly with the number of episodes  $K$ , we cannot expect to learn. Indeed, according to the lower bound (Besbes et al., 2014), the regret is necessarily linear in this case.

Table 1: Regret for optimized kernel parameters.

	bound	regret
$d = 0$	$\mathcal{R}_1$	$H^2 X \sqrt{A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}}$
	$\mathcal{R}_2$	$H^2 \sqrt{X A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}} + H^3 X^2 A \Delta^{\frac{2}{3}} K^{\frac{1}{3}}$
$d > 0$	$\mathcal{R}_1$	$H^2 \Delta^{\frac{1}{3}} K^{\frac{2d+2}{2d+3}}$
	$\mathcal{R}_2$	$H^2 \Delta^{\frac{1}{2}} K^{\frac{2d+1}{2d+2}} + H^{\frac{3}{2}} \Delta^{\frac{1}{4}} K^{\frac{3}{4}}$

variation  $\Delta$ . Recent work has started to tackle this issue in bandit algorithms (Chen et al., 2019; Auer et al., 2019), and finite MDPs using sliding windows (Cheung et al., 2020). Their extension to continuous MDPs is left to future work.

## 4 Efficient Implementation

Since **KeRNS** uses non-parametric kernel estimators, its computational complexity scales with the number of observed transitions. Let  $\tau_A$  be the time required to compute the maximum of  $a \mapsto Q_h^k(x, a)$ . Similarly to **Kernel-UCBVI**, its total space complexity is  $\mathcal{O}(KH)$  and its time complexity per episode  $k$  is  $\mathcal{O}(Hk^2 + H\tau_A k)$ , resulting in a total runtime of  $\mathcal{O}(HK^3 + H\tau_A K^2)$ . This runtime is very prohibitive in practice, especially in changing environments, where we might need to run the algorithm for a very long time. Domingues et al. (2020) propose a version of **Kernel-UCBVI** with improved per-episode time complexity of  $\mathcal{O}(H\tau_A k)$  based on real-time dynamic programming (RTDP) (Barto et al., 1995; Efroni et al., 2019). However, this requires the upper bounds  $V_h^k$  to be non-increasing, which is not the case in **KeRNS**, since  $V_h^k$  increases in regions that were not visited recently. This property is necessary to promote extra exploration and adapt to possible changes. Additionally, the RTDP-based approach of Domingues et al. (2020) still has a time complexity that scales with time, which can be a considerable issue in practice. Here, we propose an alternative to run **KeRNS** in *constant* time per episode, while controlling the impact of this speed-up on the regret.

### 4.1 Using Representative States and Actions

As proposed by Kveton and Theodorou (2012) and Barreto et al. (2016), we take an approach based on using *representative states* to construct an algorithm called **RS-KeRNS** (for **KeRNS** on Representative States). In each episode  $k$ , **RS-KeRNS** keeps and updates sets of representative states  $\bar{\mathcal{X}}_h$ , actions  $\bar{\mathcal{A}}_h$  and next-states  $\bar{\mathcal{Y}}_h$ , for each  $h$ , whose cardinalities are denoted by  $\bar{X}_h, \bar{A}_h$  and  $\bar{Y}_h$ , respectively. For simplicity, we omit the dependence on  $k$  of these sets and their cardinalities.

Every time a new transition  $\{x_h^k, a_h^k, x_{h+1}^k, r_h^k\}$  is observed, the representative sets are updated using Algorithm 2, which ensures that any two representative state-action pairs are at a distance greater than  $\varepsilon$  from each other. Similarly, it ensures that any pair of representative next-states are at a distance greater than  $\varepsilon_{\mathcal{X}}$  from each other. Then,  $(x_h^k, a_h^k)$  and  $x_{h+1}^k$  are mapped to their nearest neighbors in  $\bar{\mathcal{X}}_h \times \bar{\mathcal{A}}_h$  and  $\bar{\mathcal{Y}}_h$ , respectively, and the estimators of the rewards and transitions are updated. Consequently, we build a finite MDP, denoted by  $\bar{\mathcal{M}}_k$ , with  $\bar{\mathcal{X}}_h$  states,  $\bar{\mathcal{A}}_h$  actions and  $\bar{\mathcal{Y}}_h$  next-states, *per stage*  $h$ . The rewards and transitions of  $\bar{\mathcal{M}}_k$  can be stored in arrays of size  $\bar{X}_h \bar{A}_h$  and  $\bar{X}_h \bar{A}_h \bar{Y}_h$ , for each  $h$ .

**RS-KeRNS** is described precisely in Algorithm 4 in Appendix G. It computes a  $Q$ -function for all  $(\bar{x}, \bar{a}) \in \cup_h \bar{\mathcal{X}}_h \times \bar{\mathcal{A}}_h$  by running backward induction in  $\bar{\mathcal{M}}_k$ , which is then extended to any  $(x, a) \in \mathcal{X} \times \mathcal{A}$  by performing an interpolation step, as in **KeRNS**. In Appendix G, we explain how the rewards and transitions estimators of  $\bar{\mathcal{M}}_k$  can be updated online. Below, we provide regret and runtime guarantees for this efficient implementation.

---

### Algorithm 2 Update Representative Sets

---

- 1: **Input:**  $k, h, \bar{\mathcal{X}}_h, \bar{\mathcal{A}}_h, \bar{\mathcal{Y}}_h, \{x_h^k, a_h^k, x_{h+1}^k\}, \varepsilon, \varepsilon_{\mathcal{X}}$ .
  - 2: **if**  $\min_{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h \times \bar{\mathcal{A}}_h} \rho[(\bar{x}, \bar{a}), (x_h^k, a_h^k)] > \varepsilon$  **then**
  - 3:    $\bar{\mathcal{X}}_h = \bar{\mathcal{X}}_h \cup \{x_h^k\}$ ,  $\bar{\mathcal{A}}_h = \bar{\mathcal{A}}_h \cup \{a_h^k\}$
  - 4: **end if**
  - 5: **if**  $\min_{\bar{y} \in \bar{\mathcal{Y}}_h} \rho_{\mathcal{X}}(\bar{x}, x_{h+1}^k) > \varepsilon_{\mathcal{X}}$  **then**
  - 6:    $\bar{\mathcal{Y}}_h = \bar{\mathcal{Y}}_h \cup \{x_{h+1}^k\}$
  - 7: **end if**
- 

### 4.2 Theoretical Guarantees & Runtime

Theorem 2 states that **RS-KeRNS** enjoys the same regret bounds as **KeRNS** plus a bias term that can be controlled by  $\varepsilon$  and  $\varepsilon_{\mathcal{X}}$ , as long as we use a Gaussian kernel.

**Theorem 2** (regret of **RS-KeRNS**). *Let  $\chi_{(\eta, W)} : \mathbb{N} \rightarrow [0, 1]$ ,  $u, v \in \mathcal{X} \times \mathcal{A}$ , and consider the kernel*

$$\Gamma(t, u, v) = \chi_{(\eta, W)}(t) \exp\left(-\rho[u, v]^2 / (2\sigma^2)\right)$$

*assumed to satisfy Assumption 4. With this choice of kernel, the regret of **RS-KeRNS** is bounded by*

$$\mathcal{R}(K) \lesssim \mathcal{R}^{\text{KeRNS}}(K) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3$$

*with probability at least  $1 - \delta$ , where  $\mathcal{R}^{\text{KeRNS}}$  is regret bound of **KeRNS** given in Theorem 1.*

*Proof.* This result comes from theorems 5 and 6 in Appendix G. See Appendix B for a proof outline.

**Lemma 1** (runtime of **RS-KeRNS**). *Consider the kernel defined in Theorem 2, and let  $(\alpha_i)_{i \geq 1}$  be a sequence of real numbers in  $[0, 1]$ . If we take  $\chi_{(\eta, W)}(t) = \prod_{i=1}^t \alpha_i$ , the per-episode runtime of **RS-KeRNS** is bounded by*

$$\mathcal{O}(H \min(k^2, |\mathcal{C}_\varepsilon| |\mathcal{C}'_{\varepsilon_X}|) + H \min(k, |\mathcal{C}'_{\varepsilon_X}|) \tau_A),$$

where  $|\mathcal{C}_\varepsilon|$  is the  $\varepsilon$ -covering number of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $|\mathcal{C}'_{\varepsilon_X}|$  is the  $\varepsilon_X$ -covering number of  $(\mathcal{X}, \rho)$ , and  $\tau_A$  is the time required to compute the maximum of  $a \mapsto Q_h^k(x, a)$ . In particular, we can take  $\alpha_i = \eta$  for all  $i$ , which gives an exponential-discount strategy for handling non-stationarity.

*Proof.* By construction, in any episode  $k$ , we have  $\bar{X}_h \bar{A}_h \leq \min(k, |\mathcal{C}_\varepsilon|)$  and  $\bar{Y}_h \leq \min(k, |\mathcal{C}'_{\varepsilon_X}|)$ . Backward induction (Algorithm 5) is performed in  $\mathcal{O}(\sum_h \bar{X}_h \bar{A}_h \bar{Y}_h + \tau_A \sum_h \bar{Y}_h)$  time, and the choice of  $\chi_{(\eta, W)}(t)$  implies that the model updates can be done in  $\mathcal{O}(\sum_h \bar{X}_h \bar{A}_h \bar{Y}_h)$  time, as detailed in Appendix G.2.

Consequently, the constants  $\varepsilon$  and  $\varepsilon_X$  provide a trade-off between regret and computational complexity. Since  $|\mathcal{C}_\varepsilon| = \mathcal{O}(\varepsilon^{-d_1})$  and  $|\mathcal{C}'_{\varepsilon_X}| = \mathcal{O}(\varepsilon^{-d_2})$ , increasing  $(\varepsilon, \varepsilon_X)$  may reduce exponentially the runtime of **RS-KeRNS**, while having only a linear increase in its regret.

Kveton and Theodorou (2012) and Barreto et al. (2016) studied the idea of using representative states to accelerate kernel-based RL (KBRL), but we provide the first regret bounds in this setting. More precisely, our result improves previous work in the following aspects: (i) Kveton and Theodorou (2012) and Barreto et al. (2016) do not tackle exploration and do not have finite-time analyses: they provide approximate versions of the KBRL algorithm of Ormoneit and Sen (2002) which has asymptotic guarantees assuming that transitions are generated from independent samples; (ii) The error bounds of Kveton and Theodorou (2012) scale with  $\exp(1/\sigma^2)$ . In our online setting,  $\sigma$  can be chosen as a function of the horizon  $K$ , and their bound could result in an error that scales exponentially with  $K$ , instead of linearly, as ours. Our result comes from an improved analysis of the smoothness of kernel estimators, that leverages the regularization constant  $\beta$  (Lemma 25); (iii) Barreto et al. (2016) propose an algorithm that also builds a set of representative states in an online way. However, their theoretical guarantees only hold when this set is fixed, i.e., cannot be updated during exploration, whereas our bounds hold in this case; (iv) unlike (Kveton and Theodorou, 2012; Barreto et al., 2016), our theoretical results also hold in continuous action spaces.

### 4.3 Numerical Validation

To illustrate the behavior of **RS-KeRNS**, we consider a continuous MDP whose state-space is the unit

ball in  $\mathbb{R}^2$  with four actions, representing a move to the right, left, up or down. The agent starts at  $(0, 0)$ . Let  $b_i^k \in \{0, 0.25, 0.5, 0.75, 1\}$  and  $x_i \in \{(0.8, 0.0), (0.0, 0.8), (-0.8, 0.0), (0.0, -0.8)\}$ . We consider the following mean reward function:

$$r_h^k(x, a) = \sum_{i=1}^4 b_i^k \max\left(0, 1 - \frac{\|x - x_i\|_2}{0.5}\right)$$

which do not depend on  $h$ . Every  $N$  episodes, the coefficients  $b_i^k$  are changed, which impact the optimal policy.

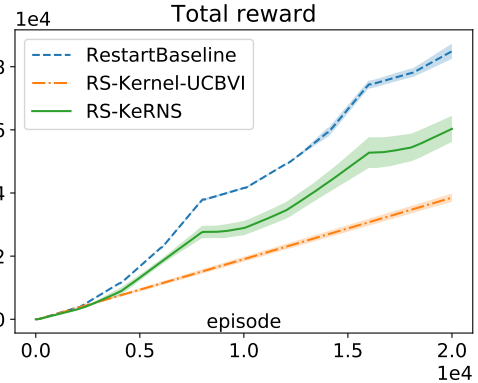


Figure 1: Performance of **RS-KeRNS** compared to baselines for  $N = 2000$ . Average over 4 runs.

Taking  $\eta = \exp(-(1/N)^{2/3})$ , we used the kernel  $\Gamma(t, u, v) = \eta^t \exp(-(\rho[u, v]/\sigma)^4/2)$ . We set  $\sigma = 0.05$ ,  $\varepsilon = \varepsilon_X = 0.1$ ,  $\beta = 0.01$ ,  $H = 15$  and ran the algorithm for  $2 \times 10^4$  episodes. **KeRNS** was compared to two baselines: (i) **Kernel-UCBVI** combined with representative states, that we call **RS-Kernel-UCBVI**, which is designed for stationary environments. This corresponds to **RS-KeRNS** with  $\chi(t) = 1$ , that is,  $\eta = 1$ ; (ii) A restart-based algorithm, called **RestartBaseline**, which is implemented as **RS-Kernel-UCBVI**, but it has full information about when the environment changes, and, at every change, it restarts its reward estimator and its bonuses. We can see that, as expected, **RS-KeRNS** outperforms **RS-Kernel-UCBVI**, which was not designed for non-stationary environments, and is able to “track” the behavior of the restart-based algorithm which has full information about how the environment changes. In Appendix I, we give more details about the experimental setup and provide more experiments, varying the period  $N$  of changes in the MDP and the kernel function.

## 5 Proof Outline

We now outline the proof of Theorem 1 assuming, for simplicity, that the rewards are known.



**Bias due to non-stationarity** To bound the bias, we introduce an average MDP with transitions  $\bar{P}_h^k$ :

$$\bar{P}_h^k(y|x, a) \stackrel{\text{def}}{=} \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) P_h^s(y|x, a) + \frac{\beta P_h^k(y|x, a)}{\mathbf{C}_h^k(x, a)},$$

where  $(P_h^s)_{s,h}$  are the true transitions at time  $(s, h)$ . We prove that, for any  $L$ -Lipschitz function  $f$  bounded by  $H$ , (Corollary 2):

$$\left| (P_h^k - \bar{P}_h^k) f(x, a) \right| \leq \mathbf{bias}_p(k, h)$$

where the term  $\mathbf{bias}_p(k, h)$  is defined as

$$L \sum_{i=1 \vee (k-W)}^{k-1} \sup_{x, a} \mathbb{W}_1(P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)) + \frac{2C_3 H}{\beta} \frac{\eta^W}{1-\eta}.$$

**Concentration** Using concentration inequalities for weighted sums, we prove that  $\hat{P}_h^k$  is close to the average transition  $\bar{P}_h^k$  using Hoeffding- and Bernstein-type inequalities (lemmas 5, 6, 7, and 8), and define an event  $\mathcal{G}$  where our confidence sets hold (Lemma 9), such that  $\mathbb{P}[\mathcal{G}] \geq 1 - \delta/2$ . For instance, Lemma 5 gives us

$$\left| (\hat{P}_h^k - \bar{P}_h^k) V_{k,h+1}^*(x, a) \right| \lesssim \sqrt{\frac{H^2}{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} + L\sigma,$$

which explains the form of the exploration bonuses.

**Upper bound on the true value function** On the event  $\mathcal{G}$ , we show that (Lemma 10):

$$Q_h^k(x, a) + \sum_{h'=h}^H \mathbf{bias}(k, h) \geq Q_{k,h}^*(x, a)$$

where the term  $\mathbf{bias}(k, h)$  is the sum of  $\mathbf{bias}_p(k, h)$  defined above, and a similar term representing the bias in the reward estimation.

**Regret bounds** Let  $(\tilde{x}_h^k, \tilde{a}_h^k)$  be the state-action pair among the previously visited ones that is the closest to  $(x_h^k, a_h^k)$ :

$$(\tilde{x}_h^k, \tilde{a}_h^k) \stackrel{\text{def}}{=} \underset{(x_h^s, a_h^s): s < k, h \in [H]}{\operatorname{argmin}} \rho[(x_h^k, a_h^k), (x_h^s, a_h^s)].$$

We show that (see proof of Lemma 11):

$$H \sum_{k=1}^K \sum_{h=1}^H \mathbb{I} \{ \rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] > 2\sigma \} \leq H^2 |\mathcal{C}_\sigma|.$$

Thus, to simplify the outline, for all  $(k, h)$ , we assume that  $\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma$  and add  $H^2 |\mathcal{C}_\sigma|$  to the

final regret bound. On the event  $\mathcal{G}$ , we prove that the regret of **KeRNS** is bounded by (lemmas 11 and 12):

$$\begin{aligned} \mathcal{R}(K) &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \\ &\quad + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH\sigma + H^2 |\mathcal{C}_\sigma| \end{aligned}$$

where we omitted factors involving  $|\mathcal{C}_\sigma|$  and  $|\mathcal{C}'_\sigma|$  (which depend on the type of bound considered,  $\mathcal{R}_1$  or  $\mathcal{R}_2$ ), and martingale terms (which are bounded by  $\approx H^{3/2} \sqrt{K}$  with probability at least  $1 - \delta/2$ ).

Using the properties of the kernel  $\Gamma$  (Assumption 4), we prove that (Lemma 16):

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} &\lesssim HK \log \frac{1}{\eta} \left( |\mathcal{C}_\sigma| + \sqrt{\frac{|\mathcal{C}_\sigma|}{\log(1/\eta)}} \right) \\ \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} &\lesssim H |\mathcal{C}_\sigma| K \log \frac{1}{\eta} \end{aligned}$$

Finally, in Corollary 5, we prove that the bias  $\sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(h, k)$  is bounded by

$$2W(\Delta^r + L\Delta^p) + \frac{2C_3(H+1)KH}{\beta} \frac{\eta^W}{1-\eta}.$$

Putting these bounds together, we prove Theorem 1. The proof of Theorem 2 is outlined in Appendix B.

## 6 Conclusion

In this paper, we introduced and analyzed **KeRNS**, the first algorithm for continuous MDPs with dynamic regret guarantees in changing environments. Building upon previous work on using representative states for kernel-based RL, we proposed **RS-KeRNS**, a practical version of **KeRNS** that runs in constant time per episode. Moreover, we provide the first analysis that quantifies the trade-off between the regret and the computational complexity of this approach. In discrete environments, our regret bound matches the existing lower bound for multi-armed bandits in terms of the number of episodes and the variation of MDP, whereas finding a lower bound in continuous environments remains an open problem.

We believe that the ideas introduced in this paper might be useful for large-scale problems. Indeed, we provide stronger online guarantees for practical kernel-based RL, which has already been shown to be empirically

successful in medium-scale environments ( $d \approx 10$ ) (Kveton and Theodorou, 2012; Barreto et al., 2016), and we show that kernel-based RL is naturally suited to tackle non-stationarity. In larger dimension, kernel-based exploration bonuses have been recently shown to enhance exploration in RL for Atari games (Badia et al., 2020), and our approach might inspire the design of bonuses for high-dimensional non-stationary environments.

### Acknowledgements

The research presented was supported by European CHIST-ERA project DELTA, French Ministry of Higher Education and Research, Nord-Pas-de-Calais Regional Council, French National Research Agency project BOLD (ANR19-CE23-0026-04), FMJH PGMO project 2018-0045, and the SFI Sachsen-Anhalt for the project RE-BCI ZS/2019/10/102024 by the Investitionsbank Sachsen-Anhalt.

### References

- Auer, P., Gajane, P., and Ortner, R. (2019). Adaptively tracking the best bandit arm with an unknown number of distribution changes. In *Conference on Learning Theory*, pages 138–158.
- Azar, M. G., Osband, I., and Munos, R. (2017). Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 263–272. JMLR.org.
- Badia, A. P., Sprechmann, P., Vitvitskyi, A., Guo, D., Piot, B., Kapturowski, S., Tieleman, O., Arjovsky, M., Pritzel, A., Bolt, A., and Blundell, C. (2020). Never give up: Learning directed exploration strategies. In *International Conference on Learning Representations*.
- Barreto, A. M., Precup, D., and Pineau, J. (2016). Practical kernel-based reinforcement learning. *The Journal of Machine Learning Research*, 17(1):2372–2441.
- Barto, A. G., Bradtko, S. J., and Singh, S. P. (1995). Learning to act using real-time dynamic programming. *Artificial intelligence*, 72(1-2):81–138.
- Besbes, O., Gur, Y., and Zeevi, A. (2014). Stochastic multi-armed-bandit problem with non-stationary rewards. In *Advances in neural information processing systems*, pages 199–207.
- Bubeck, S., Munos, R., Stoltz, G., and Szepesvári, C. (2011). X-armed bandits. *Journal of Machine Learning Research*, 12:1587–1627.
- Chen, Y., Lee, C.-W., Luo, H., and Wei, C.-Y. (2019). A new algorithm for non-stationary contextual bandits: Efficient, optimal and parameter-free. In *Conference on Learning Theory*, pages 696–726. PMLR.
- Cheung, W. C., Simchi-Levi, D., and Zhu, R. (2020). Reinforcement learning for non-stationary Markov decision processes: The blessing of (More) optimism. In *Proceedings of the 37th International Conference on Machine Learning*.
- Choi, S. P., Yeung, D.-Y., and Zhang, N. L. (2000). Hidden-mode markov decision processes for nonstationary sequential decision making. In *Sequence Learning*, pages 264–287. Springer.
- Csáji, B. C. and Monostori, L. (2008). Value function based reinforcement learning in changing markovian environments. *Journal of Machine Learning Research*, 9(Aug):1679–1709.
- Dann, C., Lattimore, T., and Brunskill, E. (2017). Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 5713–5723.
- Dick, T., Gyorgy, A., and Szepesvari, C. (2014). Online learning in markov decision processes with changing cost sequences. In *International Conference on Machine Learning*, pages 512–520. PMLR.
- Domingues, O. D., Ménard, P., Pirota, M., Kaufmann, E., and Valko, M. (2020). Regret Bounds for Kernel-Based Reinforcement Learning. *arXiv e-prints*, page arXiv:2004.05599.
- Efroni, Y., Merlis, N., Ghavamzadeh, M., and Mannor, S. (2019). Tight regret bounds for model-based reinforcement learning with greedy policies. In *Advances in Neural Information Processing Systems*, pages 12203–12213.
- Even-Dar, E., Kakade, S. M., and Mansour, Y. (2009). Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736.
- Gajane, P., Ortner, R., and Auer, P. (2018). A sliding-window algorithm for markov decision processes with arbitrarily changing rewards and transitions. *arXiv preprint arXiv:1805.10066*.
- Garivier, A. and Moulines, E. (2011). On Upper-Confidence Bound Policies For Switching Bandit Problems. In *Algorithmic Learning Theory (ALT)*, pages 174–188. PMLR.
- Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873.

- Kleinberg, R., Slivkins, A., and Upfal, E. (2019). Bandits and experts in metric spaces. *Journal of the ACM (JACM)*, 66(4):1–77.
- Kocsis, L. and Szepesvári, C. (2006). Discounted UCB. In *2nd PASCAL Challenges Workshop*.
- Kveton, B. and Theodorou, G. (2012). Kernel-based reinforcement learning on representative states. In *Twenty-Sixth AAAI Conference on Artificial Intelligence*.
- Lecarpentier, E. and Rachelson, E. (2019). Non-stationary markov decision processes, a worst-case approach using model-based reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 7214–7223.
- Li, Y. and Li, N. (2019). Online learning for markov decision processes in nonstationary environments: A dynamic regret analysis. In *2019 American Control Conference (ACC)*, pages 1232–1237. IEEE.
- Lykouris, T., Simchowitz, M., Slivkins, A., and Sun, W. (2019). Corruption robust exploration in episodic reinforcement learning. *arXiv preprint arXiv:1911.08689*.
- Neu, G., Györfy, A., Szepesvári, C., and Antos, A. (2013). Online markov decision processes under bandit feedback. *IEEE Transactions on Automatic Control*, 59(3):676–691.
- Ormonoit, D. and Sen, Š. (2002). Kernel-based reinforcement learning. *Machine Learning*, 49(2):161–178.
- Ortner, R., Gajane, P., and Auer, P. (2019). Variational regret bounds for reinforcement learning. In *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence*.
- Ortner, R. and Ryabko, D. (2012). Online regret bounds for undiscounted continuous reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 1763–1771.
- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Russac, Y., Vernade, C., and Cappé, O. (2019). Weighted linear bandits for non-stationary environments. In *Advances in Neural Information Processing Systems*, pages 12017–12026.
- Sinclair, S., Wang, T., Jain, G., Banerjee, S., and Yu, C. (2020). Adaptive discretization for model-based reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 3858–3871.
- Sinclair, S. R., Banerjee, S., and Yu, C. L. (2019). Adaptive discretization for episodic reinforcement learning in metric spaces. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(3):1–44.
- Song, Z. and Sun, W. (2019). Efficient model-free reinforcement learning in metric spaces. *arXiv preprint arXiv:1905.00475*.
- Szita, I., Takács, B., and Lőrincz, A. (2002).  $\epsilon$ -mdps: Learning in varying environments. *Journal of Machine Learning Research*, 3(Aug):145–174.
- Yu, J. Y. and Mannor, S. (2009). Online learning in markov decision processes with arbitrarily changing rewards and transitions. In *2009 international conference on game theory for networks*, pages 314–322. IEEE.
- Zanette, A. and Brunskill, E. (2019). Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *Proceedings of the 36th International Conference on Machine Learning*.

# Appendix

## Table of Contents

---

<b>A Preliminaries</b>	<b>12</b>
A.1 Notation . . . . .	12
A.2 Probabilistic model . . . . .	13
A.3 Exploration Bonuses and Kernel Backward Induction . . . . .	13
<b>B Proof Outline</b>	<b>15</b>
B.1 Theorem 2 . . . . .	15
B.2 Optimized Kernel Parameters and Regret Bounds . . . . .	16
<b>C Handling the bias due to non-stationarity</b>	<b>17</b>
<b>D Concentration</b>	<b>19</b>
D.1 Concentration inequalities for weighted sums . . . . .	19
D.2 Hoeffding-type concentration inequalities . . . . .	19
D.3 Bernstein-type concentration inequality . . . . .	23
D.4 Good event . . . . .	26
<b>E Upper bound on true value function</b>	<b>27</b>
<b>F Regret bounds</b>	<b>29</b>
F.1 Regret bound in terms of the sum of exploration bonuses (UCRL-type) . . . . .	29
F.2 Regret bound in terms of the sum of exploration bonuses (UCBVI-type) . . . . .	32
F.3 Bounding the sum of bonuses and bias . . . . .	34
F.4 Final regret bounds . . . . .	38
<b>G RS-KeRNS: An efficient version of KeRNS using representative states</b>	<b>41</b>
G.1 Definitions . . . . .	41
G.2 Online updates & runtime . . . . .	44
G.3 Regret analysis . . . . .	45
<b>H Technical Lemmas</b>	<b>57</b>
<b>I Experiments</b>	<b>61</b>
I.1 Setup . . . . .	61
I.2 Results . . . . .	61

---

## A Preliminaries

### A.1 Notation

Throughout the proof, we use the following notation when omitting constants and logarithmic terms:

**Definition 4.**

$$A \lesssim B \iff A \leq B \times \text{polynomial}(d_1, d_2, \log(K), \log(1/\delta), \beta, 1/\beta, L_r, L_p).$$

Table 2 summarizes the main notations used in the paper and in the proofs.

Table 2: Table of notations

Notation	Meaning
$\rho : (\mathcal{X} \times \mathcal{A})^2 \rightarrow \mathbb{R}_+$	metric on the state-action space $\mathcal{X} \times \mathcal{A}$
$\rho_{\mathcal{X}} : \mathcal{X}^2 \rightarrow \mathbb{R}_+$	metric on the state space $\mathcal{X}$
$\mathcal{N}(\epsilon, \mathcal{X} \times \mathcal{A}, \rho)$	$\epsilon$ -covering number of the metric space $(\mathcal{X} \times \mathcal{A}, \rho)$
$K$	number of episodes
$H$	horizon, length of each episode
$\delta$	confidence parameter
$\sigma$	kernel bandwidth parameter
$\eta$	kernel temporal decay parameter
$W$	kernel temporal window parameter
$\beta$	regularization parameter
$ \mathcal{C}_\sigma ,  \mathcal{C}'_\sigma $	$\sigma$ -covering numbers of $(\mathcal{X} \times \mathcal{A}, \rho)$ and $(\mathcal{X}, \rho_{\mathcal{X}})$ , respectively
$L_r, L_p, L$	Lipschitz constants of the rewards, transitions and value functions
$\Gamma$	kernel function from $\mathbb{N}^* \times (\mathcal{X} \times \mathcal{A})^2$ to $[0, 1]$
$\bar{\Gamma}(\eta, W)$	parameterized kernel function from $\mathbb{N}^* \times \mathbb{R}_+$ to $[0, 1]$
$C_1, C_2, C_3$	constants related to kernel properties, see Assumption 4
$G : \mathbb{R}_+ \rightarrow \mathbb{R}_+$	function related to kernel properties, see Assumption 4
$\mathcal{M}_k$	true MDP at episode $k$ , with rewards $r_h^k$ and transitions $P_h^k$
$\widehat{\mathcal{M}}_k$	empirical MDP built by <b>KeRNS</b> in episode $k$
$w_h^{k,s}(x, a)$	weight at $(x, a)$ in at time $(k, h)$ w.r.t the sample $(x_h^s, a_h^s)$ (Def. 1)
$\widehat{w}_h^{k,s}(x, a)$	normalized version of $w_h^{k,s}(x, a)$ (Def. 1)
$(Q_{k,h}^*, V_{k,h}^*)_{h \in [H]}$	true value functions in episode $k$
$(Q_h^k, V_h^k)_{h \in [H]}$	value functions computed by <b>KeRNS</b> in episode $k$
$(\widetilde{Q}_h^k, \widetilde{V}_h^k)_{h \in [H]}$	value functions computed by <b>RS-KeRNS</b> in episode $k$
$\mathcal{X}_h^k, \mathcal{A}_h^k, \mathcal{Y}_h^k$	sets of representative states, actions and next states, at stage $h$ of episode $k$
$\zeta_h^k$	mapping from $\mathcal{X} \times \mathcal{A}$ to $\mathcal{X}_h^k \times \mathcal{A}_h^k$
$\bar{\zeta}_h^k$	mapping from $\mathcal{X}$ to $\mathcal{Y}_h^k$
$\Delta^r, \Delta^p$	temporal variation of the rewards and transitions (Def. 3)
$\Delta$	temporal variation of the MDP = $\Delta^r + L\Delta^p$
$d (= d_1)$	covering dimension of $(\mathcal{X} \times \mathcal{A}, \rho)$
$d_2$	covering dimension of $(\mathcal{X}, \rho_{\mathcal{X}})$
$\epsilon$	threshold distance to add a new representative state-action pair
$\epsilon_{\mathcal{X}}$	threshold distance to add a new representative state
$\mathcal{L}(L, H)$	set of $L$ -Lipschitz functions from $\mathcal{X}$ to $\mathbb{R}$ bounded by $H$
$\mathbf{bias}_p(k, h)$	bias in transition estimation at time $(k, h)$ , (Def. 6)
$\mathbf{bias}_r(k, h)$	bias in reward estimation at time $(k, h)$ (Def. 6)
$\mathbf{bias}(k, h)$	sum of biases $\mathbf{bias}_r(k, h) + \mathbf{bias}_p(k, h)$ (Def. 6)
$\mathcal{G}$	good event, on which confidence intervals hold (Lemma 9)
$\log^+(z)$	equal to $\log(z + e)$ for any $z \in \mathbb{R}$

## A.2 Probabilistic model

The interaction between the algorithm and the MDP defines a stochastic process  $(x_h^s, a_h^s, x_{h+1}^s, \tilde{r}_h^s)$  for  $h \in [H]$  and  $s \in \mathbb{N}^*$ , representing the state, the action, the next state and the reward at step  $h$  of episode  $s$ . Let  $\mathcal{H}_h^s \stackrel{\text{def}}{=} \left\{ x_{h'}^{s'}, a_{h'}^{s'}, x_{h'+1}^{s'}, \tilde{r}_{h'}^{s'} \right\}_{s' < s, h' \in [H]} \cup \left\{ x_{h'}^s, a_{h'}^s, x_{h'+1}^s, \tilde{r}_{h'}^s \right\}_{h' < h}$  be the history of the process up to time  $(s, h)$ .

We define  $\mathcal{F}_h^s$  as the  $\sigma$ -algebra generated by  $\mathcal{H}_h^s$ , and denote its corresponding filtration by  $(\mathcal{F}_h^s)_{s,h}$ .

## A.3 Exploration Bonuses and Kernel Backward Induction

A reinforcement learning algorithm can be seen as a mapping from the set of possible histories  $\bigcup_{h \in [H], k \in \mathbb{N}^*} (\mathcal{X} \times \mathcal{A} \times \mathcal{X} \times [0, 1])^{kh-1}$  to the set of actions  $\mathcal{A}$ .

For a time  $(k, h)$ , **KeRNS** performs this mapping in the following way:

1. Build  $\hat{r}_h^k$  and  $\hat{P}_h^k$  as in Definition 2, which are  $\mathcal{F}_h^{k-1}$ -measurable.

2. For each  $h \in [H]$ , with  $V_{H+1}^k = 0$ ,

- Compute

$$\tilde{Q}_h^k(x, a) = \hat{r}_h^k(x, a) + \hat{P}_h^k V_{h+1}^k(x, a) + \mathbf{B}_h^k(x, a) \quad \text{for all } (x, a) \in \{(x_h^s, a_h^s)\}_{s < k}$$

- Define, for any  $(x, a)$ ,

$$Q_h^k(x, a) = \min_{s \in [k-1]} \left[ \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho [(x, a), (x_h^s, a_h^s)] \right]$$

$$V_h^k(x) = \min \left( H - h + 1, \max_{a'} Q_h^k(x, a') \right)$$

3. Choose the action  $a_h^k \in \operatorname{argmax}_a Q_h^k(x_h^k, a)$ .

Notice the algorithmic structure of **KeRNS** is the same as **Kernel-UCBVI** (Domingues et al., 2020). However, **KeRNS** uses non-stationary kernels to be able to adapt to changing environments.

It can be checked that Algorithm 3 returns the functions  $Q_h^k$  described above.

The exploration bonus is defined below:

**Definition 5** (exploration bonuses). *The exploration bonus in  $(x, a)$  at time  $(k, h)$  is defined as*

$$\begin{aligned} \mathbb{B}_h^k(x, a) &\stackrel{\text{def}}{=} r\mathbb{B}_h^k(x, a) + p\mathbb{B}_h^k(x, a), \quad \text{where} \\ r\mathbb{B}_h^k(x, a) &\stackrel{\text{def}}{=} \sqrt{\frac{2\Box_1^r(k, \delta/8)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_r(k, \delta/8)\sigma, \quad \text{and} \\ p\mathbb{B}_h^k(x, a) &\stackrel{\text{def}}{=} \sqrt{\frac{2H^2\Box_1^p(k, \delta/8)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_p(k, \delta/8)\sigma \end{aligned}$$

where

$$\begin{aligned} \Box_1^r(k, \delta) &= \tilde{\mathcal{O}}(d_1) = \log\left(\frac{\mathcal{N}(\sigma^2/K, \mathcal{X} \times \mathcal{A}, \rho) \sqrt{1+k/\beta}}{\delta}\right) \\ \mathbf{b}_r(k, \delta) &= \tilde{\mathcal{O}}(L + \sqrt{d_1}) = \left(\frac{C_2}{2\beta^{3/2}} \sqrt{2\Box_1^r(k, \delta)} + \frac{4C_2}{\beta}\right) + 2L_r L \left(1 + \sqrt{\log^+(C_1 k/\beta)}\right) \\ \Box_1^p(k, \delta) &= \tilde{\mathcal{O}}(d_1) = \log\left(\frac{HN(\sigma^2/KH, \mathcal{X} \times \mathcal{A}, \rho) \sqrt{1+k/\beta}}{\delta}\right) \\ \mathbf{b}_p(k, \delta) &= \tilde{\mathcal{O}}(L + \sqrt{d_1}) = \left(\frac{C_2}{2\beta^{3/2}} \sqrt{2\Box_1^p(k, \delta)} + \frac{4C_2}{\beta}\right) + 2L_p L \left(1 + \sqrt{\log^+(C_1 k/\beta)}\right). \end{aligned}$$

and where  $d_1$  is the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and, for any  $z \in \mathbb{R}$ ,  $\log^+(z) = \log(z + e)$ .

---

**Algorithm 3** Kernel Backward Induction with Exploration Bonuses

---

- 1: **Input:**  $k, H, \Gamma, L, \beta$ , transitions  $(x_h^s, a_h^s, x_{h+1}^s, \tilde{r}_h^s)_{s=1}^{k-1}$  for all  $h \in [H]$ .
  - 2: **Initialization:**  $V_{H+1}(x) = 0$  for all  $x \in \mathcal{X}$
  - 3: **for**  $h = H, \dots, 1$  **do**
  - 4:   **for**  $m = 1, \dots, k-1$  **do**
  - 5:     // Using weights given in Def.1 and bonus given in Def. 5, compute:
  - 6:      $\tilde{Q}_h(x_h^m, a_h^m) = \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x_h^m, a_h^m) (\tilde{r}_h^s + V_{h+1}(x_{h+1}^s)) + \mathbb{B}_h^k(x_h^m, a_h^m)$
  - 7:   **end for**
  - 8:   // Interpolated  $Q$ -function. Defined, but not computed, for all  $(x, a)$
  - 9:    $Q_h(x, a) = \min_{m \in [k-1]} \left( \tilde{Q}_h(x_h^m, a_h^m) + L\rho[(x, a), (x_h^m, a_h^m)] \right)$
  - 10:   **for**  $m = 1, \dots, k-1$  **do**
  - 11:      $V_h(x_h^m) = \min(H - h + 1, \max_a Q_h(x_h^m, a))$
  - 12:   **end for**
  - 13: **end for**
  - 14: **Return:**  $(Q_h)_{h \in [H]}$
-

## B Proof Outline

In this section, we outline the proof of the regret bound of **RS-KeRNS** (Theorem 2).

### B.1 Theorem 2

To prove the regret bound in Theorem 2 for **RS-KeRNS**, we consider the kernel:

$$\Gamma(t, u, v) = \chi_{(\eta, W)}(t) \phi(u, v), \text{ where } \phi(u, v) \stackrel{\text{def}}{=} \exp\left(-\rho[u, v]^2 / (2\sigma^2)\right),$$

for a given function  $\chi_{(\eta, W)} : \mathbb{N} \rightarrow [0, 1]$ . In each episode  $k$ , **RS-KeRNS** has build representative sets of states  $\bar{\mathcal{X}}_h^k$ , actions  $\bar{\mathcal{A}}_h^k$  and next states  $\bar{\mathcal{Y}}_h^k$ , for each  $h \in [H]$ . We define of the projections:

$$\zeta_h^k(x, a) \stackrel{\text{def}}{=} \underset{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k}{\text{argmin}} \rho[(x, a), (\bar{x}, \bar{a})], \quad \bar{\zeta}_h^k(y) \stackrel{\text{def}}{=} \underset{\bar{y} \in \bar{\mathcal{Y}}_h^k}{\text{argmin}} \rho_{\mathcal{X}}(y, \bar{y}).$$

from any  $(x, a, y)$  to their representatives.

Let  $\bar{W}_h^{k+1}(x, a) = \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))$ . In episode  $k+1$ , **RS-KeRNS** computes the following estimate of the rewards

$$\check{r}_h^{k+1}(x, a) = \frac{1}{\beta + \bar{W}_h^{k+1}(x, a)} \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \tilde{r}_h^s$$

and the following estimate of the transitions

$$\check{P}_h^{k+1}(y|x, a) = \frac{1}{\beta + \bar{W}_h^{k+1}(x, a)} \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \delta_{\bar{\zeta}_h^{s+1}(x_{h+1}^s)}(y).$$

which are similar to the estimates that would be computed by **KeRNS**, but using the projections  $\zeta$  and  $\bar{\zeta}$  to the representative states and actions. The values of  $\check{r}_h^{k+1}(x, a)$  and  $\check{P}_h^{k+1}(y|x, a)$  are defined for all  $(x, a, y) \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$ , but they only need to be stored for  $(x, a, y) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k \times \bar{\mathcal{Y}}_h^k$ , which corresponds to storing a *finite* representation of the MDP. The exploration bonuses of **RS-KeRNS** are defined similarly:

$$\check{B}_h^{k+1}(x, a) \stackrel{\text{def}}{=} \tilde{\mathcal{O}} \left( \frac{H}{\sqrt{\beta + \bar{W}_h^{k+1}(x, a)}} + \frac{\beta H}{\beta + \bar{W}_h^{k+1}(x, a)} + L\sigma \right)$$

We prove that the estimates used by **RS-KeRNS** are close to the ones used by **KeRNS** up to bias terms. Then, this result is used to prove that the regret bound of **RS-KeRNS** is the same as **KeRNS**, but adding a bias term multiplied by the number of episodes. For any  $(x_h^s, a_h^s)$  with  $s < k$  and  $h \in [H]$ , we show that (consequence of Lemma 18):

$$\left| \left( \hat{P}_h^k - \check{P}_h^k \right) V(x_h^s, a_h^s) \right| \lesssim L\varepsilon_{\mathcal{X}} + 8H \frac{\varepsilon}{\sigma}$$

and similar bounds are obtained for the rewards  $\check{r}_h^k(x, a)$  (Lemma 19) and the exploration bonuses (Lemma 20). This allows us to prove that the regret of **RS-KeRNS** is bounded as (theorems 5 and 6)

$$\mathcal{R}^{\text{RS-KeRNS}}(K) \lesssim \mathcal{R}^{\text{KeRNS}}(K) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3.$$

If we choose  $\chi_{(\eta, W)}(t) = \prod_{i=1}^t \alpha_i$ , the estimators used by **RS-KeRNS** can be updated online. Indeed, as detailed in Appendix G, we can related the estimates at time  $(k+1, h)$  to the ones at time  $(k, h)$ :

$$\begin{aligned} \bar{W}_h^{k+1}(\bar{x}, \bar{a}) &= \phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) + \alpha_{k-s} \bar{W}_h^k(\bar{x}, \bar{a}), \\ \check{r}_h^{k+1}(\bar{x}, \bar{a}) &= \frac{\phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) \tilde{r}_h^k}{\beta + \bar{W}_h^{k+1}(\bar{x}, \bar{a})} + \alpha_{k-s} \left( \frac{\beta + \bar{W}_h^k(\bar{x}, \bar{a})}{\beta + \bar{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{r}_h^k(\bar{x}, \bar{a}), \quad \text{and} \\ \check{P}_h^{k+1}(y|\bar{x}, \bar{a}) &= \frac{\phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) \delta_{\bar{\zeta}_h^{k+1}(x_{h+1}^k)}(y)}{\beta + \bar{W}_h^{k+1}(\bar{x}, \bar{a})} + \alpha_{k-s} \left( \frac{\beta + \bar{W}_h^k(\bar{x}, \bar{a})}{\beta + \bar{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{P}_h^k(y|\bar{x}, \bar{a}). \end{aligned}$$



One issue that we need to solve is that  $\widetilde{W}_h^k(\bar{x}, \bar{a})$ ,  $\check{r}_h^k(\bar{x}, \bar{a})$  and  $\check{P}_h^k(y|\bar{x}, \bar{a})$  were not necessarily computed before episode  $k + 1$ . This happens when  $(\bar{x}, \bar{a})$  is a new representative state-action pair added in episode  $k$ . In Section G.2, we show that this can be easily handled by defining some auxiliary quantities that can be updated online and that can be used to initialize the values  $\widetilde{W}_h^k(\bar{x}, \bar{a})$ ,  $\check{r}_h^k(\bar{x}, \bar{a})$  and  $\check{P}_h^k(y|\bar{x}, \bar{a})$  when necessary, with little overhead to the runtime of the algorithm.

## B.2 Optimized Kernel Parameters and Regret Bounds

Table 3: Regret bound for optimized kernel parameters, for  $W = \log_\eta((1 - \eta)/K)$ .

	$\sigma$	$\log\left(\frac{1}{\eta}\right)$	condition	bound	regret
$d = 0$	0	$\Delta^{\frac{2}{3}} K^{-\frac{2}{3}}$	$\Delta < K$	$\mathcal{R}_1$	$H^2 X \sqrt{A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}}$
	0	$\Delta^{\frac{2}{3}} K^{-\frac{2}{3}}$	$\Delta < K$	$\mathcal{R}_2$	$H^2 \sqrt{X A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}} + H^3 X^2 A \Delta^{\frac{2}{3}} K^{\frac{1}{3}}$
$d > 0$	$\left(\frac{1}{K}\right)^{\frac{1}{2d+3}}$	$\Delta^{\frac{2}{3}} K^{-\frac{2d+2}{2d+3}}$	$\Delta < K^{\frac{3}{2d+3}}$	$\mathcal{R}_1$	$H^2 \Delta^{\frac{1}{3}} K^{\frac{2d+2}{2d+3}}$
	$\left(\frac{1}{K}\right)^{\frac{1}{2d+2}}$	$\frac{\Delta^{\frac{1}{2}}}{H} K^{-\frac{2d+1}{2d+2}}$	$\Delta < K^{\frac{1}{d+1}}$	$\mathcal{R}_2$	$H^2 \Delta^{\frac{1}{2}} K^{\frac{2d+1}{2d+2}} + H^{\frac{3}{2}} \Delta^{\frac{1}{4}} K^{\frac{3}{4}}$

## C Handling the bias due to non-stationarity

**Lemma 2** (temporal bias). *Let  $(F_h^s)_{h,s}$  be an arbitrary sequence of functions from  $\mathcal{X} \times \mathcal{A}$  to  $\mathbb{R}$  bounded by  $M$ . Then,*

$$\left| \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) (F_h^s(x, a) - F_h^k(x, a)) \right| \leq \sum_{i=1 \vee (k-W)}^{k-1} |F_h^i(x, a) - F_h^{i+1}(x, a)| + \frac{2MC_3}{\beta} \frac{\eta^W}{1-\eta}.$$

*Proof.* The result is straightforward when  $k \leq W$ . Assuming  $k > W$ , we have

$$\begin{aligned} & \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) (F_h^s(x, a) - F_h^k(x, a)) \\ &= \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=k-W}^{k-1} w_h^{k,s}(x, a) (F_h^s(x, a) - F_h^k(x, a)) + \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-W-1} w_h^{k,s}(x, a) (F_h^s(x, a) - F_h^k(x, a)) \\ &\leq \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=k-W}^{k-1} w_h^{k,s}(x, a) \sum_{i=s}^{k-1} (F_h^i(x, a) - F_h^{i+1}(x, a)) + \frac{2MC_3}{\beta} \sum_{s=1}^{k-W-1} \eta^{k-1-s} \\ &\leq \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{i=k-W}^{k-1} \left( \sum_{s=k-W}^i w_h^{k,s}(x, a) \right) (F_h^i(x, a) - F_h^{i+1}(x, a)) + \frac{2MC_3}{\beta} \frac{\eta^W - \eta^{k-1}}{1-\eta} \\ &\leq \sum_{i=k-W}^{k-1} |F_h^i(x, a) - F_h^{i+1}(x, a)| + \frac{2MC_3}{\beta} \frac{\eta^W}{1-\eta}, \end{aligned}$$

where in the first inequality we used by Assumption 4 that

$$\begin{aligned} w_h^{k,s}(x, a) &= \Gamma(k-s-1, (x, a), (x_h^s, a_h^s)) \\ &= \bar{\Gamma}_{(\eta, W)}(k-s-1, \rho[(x, a), (x_h^s, a_h^s)]) \\ &\leq C_3 \eta^{k-s-1}. \end{aligned}$$

By symmetry, we obtain

$$\frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) (F_h^k(x, a) - F_h^s(x, a)) \leq \sum_{i=k-W}^{k-1} |F_h^i(x, a) - F_h^{i+1}(x, a)| + \frac{2MC_3}{\beta} \frac{\eta^W}{1-\eta}.$$

which concludes the proof.  $\square$

**Definition 6** (temporal bias of the MDP). *The temporal bias at time  $(k, h)$  is defined by*

$$\mathbf{bias}(k, h) = \mathbf{bias}_r(k, h) + \mathbf{bias}_p(k, h)$$

where

$$\begin{aligned} \mathbf{bias}_r(k, h) &= \sum_{i=1 \vee (k-W)}^{k-1} \sup_{x, a} |r_h^i(x, a) - r_h^{i+1}(x, a)| + \frac{2C_3}{\beta} \frac{\eta^W}{1-\eta} \\ \mathbf{bias}_p(k, h) &= L \sum_{i=1 \vee (k-W)}^{k-1} \sup_{x, a} \mathbb{W}_1(\mathbb{P}_h^i(\cdot|x, a), \mathbb{P}_h^{i+1}(\cdot|x, a)) + \frac{2C_3 H}{\beta} \frac{\eta^W}{1-\eta} \end{aligned}$$

**Definition 7** (Average MDP at episode  $k$ ). *Let*

$$\bar{r}_h^k(x, a) \stackrel{\text{def}}{=} \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) r_h^s(x, a) + \frac{\beta}{\mathbf{C}_h^k(x, a)} r_h^k(y|x, a)$$

$$\bar{P}_h^k(y|x, a) \stackrel{\text{def}}{=} \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) P_h^s(y|x, a) + \frac{\beta}{\mathbf{C}_h^k(x, a)} P_h^k(y|x, a)$$

and let  $\mathcal{M}_k^{\text{av}}$  be the MDP with transitions  $\{\bar{P}_h^k\}_h$  and rewards  $\{\bar{r}_h^k\}_h$ .

**Corollary 2.** *Let  $\mathcal{L}(L, H)$  be the class of  $L$ -Lipschitz functions from  $\mathcal{X}$  to  $\mathbb{R}$  bounded by  $H$ . Then,*

$$\sup_{x, a} |r_h^k(x, a) - \bar{r}_h^k(x, a)| \leq \mathbf{bias}_r(k, h)$$

$$\sup_{f \in \mathcal{L}(L, H)} \left| \left( P_h^k - \bar{P}_h^k \right) f(x, a) \right| \leq \mathbf{bias}_p(k, h)$$

*Proof.* For the reward term, we have from Lemma 2:

$$\begin{aligned} |r_h^k(x, a) - \bar{r}_h^k(x, a)| &= \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) (r_h^s(x, a) - \bar{r}_h^s(x, a)) \right| \\ &\leq \sum_{i=1 \vee (k-W)}^{k-1} |r_h^i(x, a) - r_h^{i+1}(x, a)| + \frac{2C_3}{\beta} \frac{\eta^W}{1-\eta} \\ &\leq \mathbf{bias}_r(k, h). \end{aligned}$$

For the transitions term, we also apply Lemma 2 and the definition of the 1-Wasserstein distance:

$$\begin{aligned} \left| \left( P_h^k - \bar{P}_h^k \right) f(x, a) \right| &= \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) \left( P_h^s f(x, a) - \bar{P}_h^s f(x, a) \right) \right| \\ &\leq \sum_{i=1 \vee (k-W)}^{k-1} |P_h^i f(x, a) - P_h^{i+1} f(x, a)| + \frac{2C_3 H}{\beta} \frac{\eta^W}{1-\eta} \\ &\leq L \sum_{i=1 \vee (k-W)}^{k-1} \mathbb{W}_1(P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)) + \frac{2C_3 H}{\beta} \frac{\eta^W}{1-\eta} \\ &\leq \mathbf{bias}_p(k, h). \end{aligned}$$

**Remark 1.** *Since the functions in  $\mathcal{L}(L, H)$  are bounded, the 1-Wasserstein distance could be replaced by the total variation (TV) distance  $\|P_h^i(\cdot|x, a) - P_h^{i+1}(\cdot|x, a)\|_1$ .*

□

## D Concentration

In this Section, we provide confidence intervals that will be used to prove our regret bounds. The main concentration results are presented in Lemma 9, which defines an event  $\mathcal{G}$  where all the confidence intervals hold, and we show that  $\mathbb{P}[\mathcal{G}] \geq 1 - \delta/2$ .

### D.1 Concentration inequalities for weighted sums

We reproduce here the concentration inequalities for weighted sums proved by Domingues et al. (2020), which we will need.

**Lemma 3** (Hoeffding type inequality (Domingues et al., 2020)). *Consider the sequences of random variables  $(w_t)_{t \in \mathbb{N}^*}$  and  $(Y_t)_{t \in \mathbb{N}^*}$  adapted to a filtration  $(\mathcal{F}_t)_{t \in \mathbb{N}}$ . Assume that, for all  $t \geq 1$ ,  $w_t$  is  $\mathcal{F}_{t-1}$  measurable and  $\mathbb{E} \left[ \exp(\lambda Y_t) \middle| \mathcal{F}_{t-1} \right] \leq \exp(\lambda^2 c^2 / 2)$  for all  $\lambda > 0$ .*

Let  $S_t \stackrel{\text{def}}{=} \sum_{s=1}^t w_s Y_s$  and  $V_t \stackrel{\text{def}}{=} \sum_{s=1}^t w_s^2$ , and assume  $w_s \leq 1$  almost surely for all  $s$ . Then, for any  $\beta > 0$ , with probability at least  $1 - \delta$ , for all  $t \geq 1$ ,

$$\frac{|S_t|}{\sum_{s=1}^t w_s + \beta} \leq \sqrt{2c^2 \log \left( \frac{\sqrt{1+t/\beta}}{\delta} \right) \frac{1}{\sum_{s=1}^t w_s + \beta}}.$$

*Proof.* See Lemma 2 of Domingues et al. (2020). □

**Lemma 4** (Bernstein type inequality (Domingues et al., 2020)). *Consider the sequences of random variables  $(w_t)_{t \in \mathbb{N}^*}$  and  $(Y_t)_{t \in \mathbb{N}^*}$  adapted to a filtration  $(\mathcal{F}_t)_{t \in \mathbb{N}}$ . Let*

$$S_t \stackrel{\text{def}}{=} \sum_{s=1}^t w_s Y_s, \quad V_t \stackrel{\text{def}}{=} \sum_{s=1}^t w_s^2 \mathbb{E} \left[ Y_s^2 \middle| \mathcal{F}_{s-1} \right] \quad \text{and} \quad W_t \stackrel{\text{def}}{=} \sum_{s=1}^t w_s,$$

Assume that, for all  $t \geq 1$ , (i)  $w_t$  is  $\mathcal{F}_{t-1}$  measurable, (ii)  $\mathbb{E} \left[ Y_t \middle| \mathcal{F}_{t-1} \right] = 0$ , (iii)  $w_t \in [0, 1]$  almost surely, (iv) there exists  $b > 0$  such that  $|Y_t| \leq b$  almost surely. Then, for all  $\beta > 0$ , with probability at least  $1 - \delta$ , for all  $t \geq 1$ ,

$$\frac{|S_t|}{\beta + \sum_{s=1}^t w_s} \leq \sqrt{2 \log(4e(2t+1)/\delta) \frac{V_t + b^2}{\left(\beta + \sum_{s=1}^t w_s\right)^2}} + \frac{2b \log(4e(2t+1)/\delta)}{3 \frac{\beta + \sum_{s=1}^t w_s}{\beta + \sum_{s=1}^t w_s}}.$$

*Proof.* See Lemma 3 of Domingues et al. (2020). □

### D.2 Hoeffding-type concentration inequalities

**Lemma 5.** *For all  $(x, a, k, h) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$ , we have*

$$\left| (\widehat{P}_h^k - \overline{P}_h^k) V_{k,h+1}^*(x, a) \right| \leq \sqrt{\frac{2H^2 \square_1^p(k, \delta)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_p(k, \delta) \sigma$$

with probability at least  $1 - \delta$ , where

$$\square_1^p(k, \delta) = \tilde{\mathcal{O}}(d_1) = \log \left( \frac{KHN(\sigma^2/(KH), \mathcal{X} \times \mathcal{A}, \rho) \sqrt{1+k/\beta}}{\delta} \right)$$

$$\mathbf{b}_p(k, \delta) = \tilde{\mathcal{O}}(L + \sqrt{d_1}) = \left( \frac{C_2}{2\beta^{3/2}} \sqrt{2\square_1^p(k, \delta)} + \frac{4C_2}{\beta} \right) + 2L_p L \left( 1 + \sqrt{\log(C_1 k/\beta)} \right)$$

and where  $d_1$  is the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ .

*Proof.* Let  $V = V_{k,h+1}^*$ . For fixed  $(x, a, h)$ , we have

$$\begin{aligned}
 & \left| (\widehat{P}_h^k - \overline{P}_h^k) V_{k,h+1}^*(x, a) \right| \\
 &= \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( V(x_{h+1}^s) - \int_{\mathcal{X}} V(y) dP_h^s(y|x, a) \right) - \frac{\beta}{\mathbf{C}_h^k(x, a)} \int_{\mathcal{X}} V(y) dP_h^k(y|x, a) \right| \\
 &\leq \underbrace{\left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( V(x_{h+1}^s) - \int_{\mathcal{X}} V(y) dP_h^s(y|x_h^s, a_h^s) \right) \right|}_{\textcircled{1}} \\
 &\quad + \underbrace{\left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( \int_{\mathcal{X}} V(y) dP_h^s(y|x_h^s, a_h^s) - \int_{\mathcal{X}} V(y) dP_h^s(y|x, a) \right) \right|}_{\textcircled{2}} \\
 &\quad + \frac{\beta H}{\mathbf{C}_h^k(x, a)}.
 \end{aligned}$$

**Bounding ① (martingale term)** Let  $Y_s = V(x_{h+1}^s) - P_h^s V(x_h^s, a_h^s)$ . From Lemma 3, we have, for a fixed tuple  $(x, a, k, h)$ ,

$$\textcircled{1} = \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) Y_s \right| \leq \sqrt{2H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right) \frac{1}{\mathbf{C}_h^k(x, a)}}$$

with probability at least  $1 - \delta$ , since  $(Y_s)_s$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^s)_s$ .

From Lemma 24, we verify that the functions

$$(x, a) \mapsto \sqrt{1/\mathbf{C}_h^k(x, a)} \quad \text{and} \quad (x, a) \mapsto \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) Y_s$$

are Lipschitz continuous, with Lipschitz constants bounded by  $C_2 k / (2\sigma\beta^{3/2})$  and  $4HC_2 k / (\beta\sigma)$ , respectively. Let  $\mathcal{C}_{\mathcal{X} \times \mathcal{A}}(\sigma^2/KH)$  be a  $(\sigma^2/KH)$ -covering of  $\mathcal{X} \times \mathcal{A}$ . Using the Lipschitz continuity of the functions above and a union bound over  $\mathcal{C}_{\mathcal{X} \times \mathcal{A}}(\sigma^2/(KH))$  and over  $h \in [H]$ , we have

$$\begin{aligned}
 \textcircled{1} &= \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) Y_s \right| \leq \sqrt{2H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right) \frac{1}{\mathbf{C}_h^k(x, a)}} \\
 &\quad + \left( \frac{C_2 k}{2\sigma\beta^{3/2}} \sqrt{2H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right)} + \frac{4HC_2 k}{\beta\sigma} \right) \frac{\sigma^2}{KH}
 \end{aligned}$$

for all  $(x, a, k, h)$  with probability at least  $1 - \delta KHN(\sigma^2/(KH), \mathcal{X} \times \mathcal{A}, \rho)$ .

**Bounding ② (spatial bias term)** We have

$$\begin{aligned}
 \textcircled{2} &= \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( \int_{\mathcal{X}} V(y) dP_h^s(y|x_h^s, a_h^s) - \int_{\mathcal{X}} V(y) dP_h^s(y|x, a) \right) \right| \\
 &\leq L \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \mathbb{W}_1(P_h^s(\cdot|x_h^s, a_h^s), P_h^s(\cdot|x, a)) \quad \text{by the definition of } \mathbb{W}_1(\cdot, \cdot) \\
 &\leq L_p L \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \rho[(x_h^s, a_h^s), (x, a)] \quad \text{by Assumption 2} \\
 &\leq 2\sigma L_p L \left( 1 + \sqrt{\log^+(C_1 k/\beta)} \right) \quad \text{by Lemma 23.}
 \end{aligned}$$

Putting together the bounds for ① and ② concludes the proof.  $\square$

**Lemma 6.** For all  $(x, a, k, h) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$ , we have

$$|\widehat{r}_h^k(x, a) - \bar{r}_h^k(x, a)| \leq \sqrt{\frac{2\Box_1^r(k, \delta)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_r(k, \delta)\sigma$$

with probability at least  $1 - \delta$ , where

$$\begin{aligned} \Box_1^r(k, \delta) &= \tilde{\mathcal{O}}(d_1) = \log \left( \frac{\mathcal{N}(\sigma^2/K, \mathcal{X} \times \mathcal{A}, \rho) \sqrt{1+k/\beta}}{\delta} \right) \\ \mathbf{b}_r(k, \delta) &= \tilde{\mathcal{O}}(L + \sqrt{d_1}) = \left( \frac{C_2}{2\beta^{3/2}} \sqrt{2\Box_1^r(k, \delta)} + \frac{4C_2}{\beta} \right) + 2L_r L \left( 1 + \sqrt{\log(C_1 k/\beta)} \right) \end{aligned}$$

and where  $d_1$  is the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ .

*Proof.* Almost identical to the proof of Lemma 5, except for the fact that the rewards are bounded by 1 instead of  $H$ .  $\square$

**Lemma 7.** Let  $\mathcal{L}(2L, 2H)$  be the class of  $2L$ -Lipschitz functions from  $\mathcal{X}$  to  $\mathbb{R}$  bounded by  $2H$ . With probability at least  $1 - \delta$ , for all  $(x, a, h, k) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$  and for all  $f \in \mathcal{L}(2L, 2H)$ , we have

$$\left| (\widehat{P}_h^k - \bar{P}_h^k) f(x, a) \right| \leq \sqrt{\frac{8H^2 \Box_2^p(k, \delta)}{\mathbf{C}_h^k(x, a)}} + \frac{2\beta H}{\mathbf{C}_h^k(x, a)} + \theta_b^1(k, \delta)\sigma^{1+d_2/2} + \theta_b^2(k, \delta)\sigma$$

where

$$\begin{aligned} \Box_2^p(k, \delta) &= \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1 d_2) = \log \left( \frac{KH \mathcal{N}(\sigma^{2+d_2/2}/KH, \mathcal{X} \times \mathcal{A}, \rho) \sqrt{1+k/\beta}}{\delta} \left( \frac{2H}{L\sigma} \right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho_{\mathcal{X}})} \right) \\ \theta_b^1(k, \delta) &= \tilde{\mathcal{O}}(\sqrt{|\mathcal{C}'_\sigma|} + \sqrt{d_1 d_2}) = \frac{4C_2}{\beta} + \frac{C_2}{2\beta^{3/2}} \sqrt{8\Box_2^p(k, \delta)} \\ \theta_b^2(k, \delta) &= \tilde{\mathcal{O}}(L) = 4L_p L \left( 1 + \sqrt{\log^+(C_1 k/\beta)} \right) + 32L \end{aligned}$$

and where  $d_1$  is the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $d_2$  is the covering dimension of  $(\mathcal{X}, \rho_{\mathcal{X}})$  and  $|\mathcal{C}'_\sigma| = \mathcal{N}(\sigma, \mathcal{X}, \rho_{\mathcal{X}})$ .

*Proof.* Fix a function  $f \in \mathcal{L}(2L, 2H)$ . Proceeding as in the proof of Lemma 5, we have

$$\begin{aligned} &\left| (\widehat{P}_h^k - \bar{P}_h^k) f(x, a) \right| \\ &= \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) \left( f(x_{h+1}^s) - \int_{\mathcal{X}} f(y) d\mathbf{P}_h^s(y|x, a) \right) - \frac{\beta}{\mathbf{C}_h^k(x, a)} \int_{\mathcal{X}} f(y) d\mathbf{P}_h^k(y|x, a) \right| \\ &\leq \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| + 4\sigma L_p L \left( 1 + \sqrt{\log^+(C_1 k/\beta)} \right) + \frac{2\beta H}{\mathbf{C}_h^k(x, a)}. \end{aligned}$$

where  $Y_s(f) = f(x_{h+1}^s) - \mathbf{P}_h^s f(x_h^s, a_h^s)$ .

Now, for fixed  $(x, a, k, h, f)$ , Lemma 3 gives us

$$\left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| \leq \sqrt{8H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right) \frac{1}{\mathbf{C}_h^k(x, a)}}$$

with probability at least  $1 - \delta$ , since  $(Y_s(f))_s$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^s)_s$  bounded by  $4H$ .

**Covering  $\mathcal{L}(2L, 2H)$**  Now let  $\mathcal{C}_{\mathcal{L}}$  be a  $8L\sigma$ -covering of  $(\mathcal{L}(2L, 2H), \|\cdot\|_{\infty})$ . Using the fact that the function  $f \mapsto Y_s(f)$  is 2-Lipschitz with respect to  $\|\cdot\|_{\infty}$ , we do a union bound over  $\mathcal{C}_{\mathcal{L}}$  to obtain

$$\left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| \leq \sqrt{8H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right) \frac{1}{\mathbf{C}_h^k(x, a)}} + 32L\sigma$$

for all  $k$  and all  $f \in \mathcal{L}(2L, 2H)$ , with probability at least  $1 - \delta \left(\frac{2H}{L\sigma}\right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho_X)}$ , since the  $8L\sigma$ -covering number of  $\mathcal{C}_{\mathcal{L}}$  is bounded by  $\left(\frac{2H}{L\sigma}\right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho_X)}$ , by Lemma 5 of Domingues et al. (2020).

**Covering  $(\mathcal{X} \times \mathcal{A}, \rho)$**  By Lemma 24, the functions

$$(x, a) \mapsto \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| \quad \text{and} \quad (x, a) \mapsto \sqrt{\frac{1}{\mathbf{C}_h^k(x, a)}}$$

are  $4HC_2k/(\beta\sigma)$ -Lipschitz and  $C_2k/(2\beta^{3/2}\sigma)$ , respectively, with respect to the distance  $\rho$ . Let  $\mathcal{C}_{\mathcal{X} \times \mathcal{A}}$  be a  $\sigma^{2+d_2/2}/KH$  covering of  $(\mathcal{X} \times \mathcal{A}, \rho)$ . Using the continuity of the functions above, a union bound over  $\mathcal{C}_{\mathcal{X} \times \mathcal{A}}$  gives us<sup>9</sup>

$$\begin{aligned} \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| &\leq \sqrt{8H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right) \frac{1}{\mathbf{C}_h^k(x, a)}} + 32L\sigma \\ &\quad + \frac{\sigma^{2+d_2/2}}{KH} \left( \frac{4HC_2k}{\beta\sigma} + \frac{C_2k}{2\beta^{3/2}\sigma} \sqrt{8H^2 \log \left( \frac{\sqrt{1+k/\beta}}{\delta} \right)} \right) \end{aligned}$$

for all  $k$ , all  $f \in \mathcal{L}(2L, 2H)$  and all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ , with probability at least

$$1 - \delta \left(\frac{2H}{L\sigma}\right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho_X)} \mathcal{N}\left(\sigma^{2+d_2/2}/KH, \mathcal{X} \times \mathcal{A}, \rho\right)$$

and a union bound over  $(k, h) \in [K] \times [H]$  concludes the proof. □

---

<sup>9</sup>see, for instance, Lemma 6 of Domingues et al. (2020).

### D.3 Bernstein-type concentration inequality

**Lemma 8.** Let  $\mathcal{L}(2L, 2H)$  be the class of  $2L$ -Lipschitz functions from  $\mathcal{X}$  to  $\mathbb{R}$  bounded by  $2H$ . With probability at least  $1 - \delta$ , for all  $(x, a, h, k) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$  and for all  $f \in \mathcal{L}(2L, 2H)$ , we have

$$\begin{aligned} \left| (\widehat{P}_h^k - \overline{P}_h^k) f(x, a) \right| &\leq \frac{1}{H} \mathbb{P}_h^k |f|(x, a) + \frac{14H^2 C_2 \square_3(k, \delta) + 2\beta H}{\mathbf{C}_h^k(x, a)} \\ &\quad + \theta_b^3(k, \delta) \sigma^{1+d_2} + \theta_b^4(k, \delta) \sigma + \frac{2}{H} \mathbf{bias}_p(k, h) \end{aligned}$$

where  $d_1$  is the covering dimension of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $d_2$  is the covering dimension of  $(\mathcal{X}, \rho_{\mathcal{X}})$  and

$$\begin{aligned} \square_3(k, \delta) &= \widetilde{\mathcal{O}}(|\mathcal{C}'_{\sigma}| + d_1 d_2) = \log \left( \frac{4e(2k+1)}{\delta} K H N \left( \frac{\sigma^{2+d_2}}{H^2 K}, \mathcal{X} \times \mathcal{A}, \rho \right) \left( \frac{2H}{L\sigma} \right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho_{\mathcal{X}})} \right) \\ \theta_b^3(k, \delta) &= \widetilde{\mathcal{O}}(|\mathcal{C}'_{\sigma}| + d_1 d_2 + L\sigma) = \frac{2L_p L \sigma}{H^2 K} + \frac{4C_2}{H\beta} + \frac{14 \square_3(k, \delta) C_2}{\beta^2} \\ \theta_b^4(k, \delta) &= \widetilde{\mathcal{O}}(L) = 32L + 6L_p L \left( 1 + \sqrt{\log^+(C_1 k / \beta)} \right) \end{aligned}$$

where  $|\mathcal{C}'_{\sigma}| = \mathcal{O}(1/\sigma^{d_2})$  is the  $\sigma$ -covering number of  $(\mathcal{X}, \rho_{\mathcal{X}})$ .

*Proof.* We have

$$\begin{aligned} &\left| (\widehat{P}_h^k - \overline{P}_h^k) f(x, a) \right| \\ &= \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) (f(x_{h+1}^s) - \mathbb{P}_h^s f(x, a)) - \frac{\beta \mathbb{P}_h^k f(x, a)}{\mathbf{C}_h^k(x, a)} \right| \\ &\leq \underbrace{\left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( f(x_{h+1}^s) - \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x_h^s, a_h^s) \right) \right|}_{\textcircled{1}} \\ &\quad + \underbrace{\left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x_h^s, a_h^s) - \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x, a) \right) \right|}_{\textcircled{2}} + \frac{2\beta H}{\mathbf{C}_h^k(x, a)}. \end{aligned}$$

**Bounding ② (spatial bias term)** As in the proof of Lemma 5, we can show that

$$\textcircled{2} = \left| \sum_{s=1}^{k-1} \widetilde{w}_h^{k,s}(x, a) \left( \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x_h^s, a_h^s) - \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x, a) \right) \right| \leq 4\sigma L_p L \left( 1 + \sqrt{\log^+(C_1 k / \beta)} \right)$$

**Bounding the martingale term ① with a Bernstein-type inequality** Notice that  $(x, a) \mapsto \int_{\mathcal{X}} f(y) d\mathbb{P}_h^k(y|x, a)$  is bounded by  $2H$  and

$$\mathbb{E} [f(x_{h+1}^s) | \mathcal{F}_h^s] = \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x_h^s, a_h^s).$$

The conditional variance of  $f(x_{h+1}^s)$  is bounded as follows

$$\begin{aligned} \mathbb{V} [f(x_{h+1}^s) | \mathcal{F}_h^s] &= \mathbb{E} [f(x_{h+1}^s)^2 | \mathcal{F}_h^s] - \left( \int_{\mathcal{X}} f(y) d\mathbb{P}_h^s(y|x_h^s, a_h^s) \right)^2 \\ &\leq 2H \mathbb{E} [|f(x_{h+1}^s)| | \mathcal{F}_h^s] \\ &= 2H \int_{\mathcal{X}} |f(y)| d\mathbb{P}_h^s(y|x_h^s, a_h^s) \end{aligned}$$



which we use to bound its weighted average

$$\begin{aligned}
 & \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a)^2 \mathbb{V} [f(x_{h+1}^s) | \mathcal{F}_h^s] \\
 & \leq \frac{1}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) \mathbb{V} [f(x_{h+1}^s) | \mathcal{F}_h^s] \\
 & \leq \frac{2H}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) \int_{\mathcal{X}} |f(y)| d\mathbb{P}_h^s(y | x_h^s, a_h^s) \\
 & = \frac{2H}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) \mathbb{P}_h^s |f|(x, a) + \frac{2H}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) (\mathbb{P}_h^s |f|(x_h^s, a_h^s) - \mathbb{P}_h^s |f|(x, a)) \\
 & \leq 2H \left( \bar{\mathbb{P}}_h^k |f|(x, a) - \frac{\beta \mathbb{P}_h^k |f|(x, a)}{\mathbf{C}_h^k(x, a)} \right) + \frac{4HL_p L}{\mathbf{C}_h^k(x, a)} \sum_{s=1}^{k-1} w_h^{k,s}(x, a) \rho[(x_h^s, a_h^s), (x, a)] \\
 & \leq 2H \bar{\mathbb{P}}_h^k |f|(x, a) + 8HL_p L \sigma \left( 1 + \sqrt{\log^+(C_1 k / \beta)} \right)
 \end{aligned}$$

where, in the last inequality, we used Lemma 23.

Let  $\Delta(k, \delta) = \log(4e(2k+1)/\delta)$ . Let  $Y_s(f) = f(x_{h+1}^s) - \mathbb{P}_h^s f(x_h^s, a_h^s)$ . By Lemma 4, we have

$$\textcircled{1} = \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| \leq \sqrt{2\Delta(k, \delta) \frac{\sum_{s=1}^{k-1} w_h^{k,s}(x, a)^2 \mathbb{V} [f(x_{h+1}^s) | \mathcal{F}_h^s]}{\mathbf{C}_h^k(x, a)^2}} + \frac{10H\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)}$$

with probability at least  $1 - \delta$ , since, for a fixed  $f$ ,  $(Y_s(f))_s$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^s)_s$ . Using the fact that  $\sqrt{uv} \leq (u+v)/2$  for all  $u, v > 0$ ,

$$\begin{aligned}
 \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| & \leq \frac{4H^2 \Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} + \frac{1}{4H^2} \frac{\sum_{s=1}^{k-1} w_h^{k,s}(x, a)^2 \mathbb{V} [f(x_{h+1}^s) | \mathcal{F}_h^s]}{\mathbf{C}_h^k(x, a)} + \frac{10H\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} \\
 & \leq \frac{1}{H} \int_{\mathcal{X}} |f(y)| d\bar{\mathbb{P}}_h^k(y | x, a) + \frac{(4H^2 + 10H)\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} + \frac{2L_p L \sigma}{H} \left( 1 + \sqrt{\log^+(C_1 t / \beta)} \right)
 \end{aligned}$$

From Corollary 2, we have

$$\int_{\mathcal{X}} |f(y)| d\bar{\mathbb{P}}_h^k(y | x, a) = (\bar{\mathbb{P}}_h^k - \mathbb{P}_h^k) |f(y)|(x, a) + \mathbb{P}_h^k |f(y)|(x, a) \leq 2 \mathbf{bias}_p(k, h) + \mathbb{P}_h^k |f(y)|(x, a)$$

which gives us

$$\begin{aligned}
 \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| & \leq \frac{1}{H} \mathbb{P}_h^k |f(y)|(x, a) + \frac{(4H^2 + 10H)\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} \\
 & \quad + \frac{2}{H} \mathbf{bias}_p(k, h) + \frac{2L_p L \sigma}{H} \left( 1 + \sqrt{\log^+(C_1 t / \beta)} \right)
 \end{aligned}$$

with probability  $1 - \delta$ .

**Covering of  $\mathcal{X} \times \mathcal{A}$**  As a consequence of Assumption 2, the function  $(x, a) \mapsto (1/H)\mathbb{P}_h^k |f(y)|(x, a)$  is  $2L_p L$ -Lipschitz. Also, the functions

$$(x, a) \mapsto \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| \quad \text{and} \quad (x, a) \mapsto \frac{1}{\mathbf{C}_h^k(x, a)}$$

are  $4HC_2k/(\beta\sigma)$ -Lipschitz and  $C_2k/(\beta^2\sigma)$ , respectively, by Lemma 24. Consequently, a union bound over a  $(\sigma^{2+d_2}/(H^2K))$ -covering of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and over  $[H]$  gives us

$$\begin{aligned} \left| \sum_{s=1}^{k-1} \tilde{w}_h^{k,s}(x, a) Y_s(f) \right| &\leq \frac{1}{H} \mathbb{P}_h^k |f(y)|(x, a) + \frac{(4H^2 + 10H)\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} \\ &\quad + \frac{2}{H} \mathbf{bias}_p(k, h) + \frac{2L_p L \sigma}{H} \left( 1 + \sqrt{\log^+(C_1 t / \beta)} \right) \\ &\quad + \left( 2L_p L + \frac{4HC_2k}{\beta\sigma} + \frac{(4H^2 + 10H)\Delta(k, \delta)C_2k}{\beta^2\sigma} \right) \frac{\sigma^{2+d_2}}{H^2K} \end{aligned}$$

for all  $(x, a, h, k)$  with probability at least  $1 - \delta K H N \left( \frac{\sigma^{2+d_2}}{H^2K}, \mathcal{X} \times \mathcal{A}, \rho \right)$ .

**Covering of  $\mathcal{L}(2L, 2H)$**  The bounds for ① and ② give us

$$\begin{aligned} \left| (\hat{P}_h^k - \bar{P}_h^k) f(x, a) \right| &\leq \frac{1}{H} \mathbb{P}_h^k |f(y)|(x, a) + \frac{(4H^2 + 10H)\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} \\ &\quad + \frac{2}{H} \mathbf{bias}_p(k, h) + \left( 2L_p L + \frac{4HC_2k}{\beta\sigma} + \frac{(4H^2 + 10H)\Delta(k, \delta)C_2k}{\beta^2\sigma} \right) \frac{\sigma^{2+d_2}}{H^2K} \\ &\quad + 6\sigma L_p L \left( 1 + \sqrt{\log^+(C_1 k / \beta)} \right) + \frac{2\beta H}{\mathbf{C}_h^k(x, a)}. \end{aligned}$$

The  $8L\sigma$ -covering number of  $\mathcal{L}(2L, 2H)$  with respect to the infinity norm is bounded by  $(2H/(L\sigma))^{\mathcal{N}(\sigma, \mathcal{X}, \rho, x)}$ , by Lemma 5 of Domingues et al. (2020). The functions  $f \mapsto |(\mathbb{P}_h^k - \hat{P}_h^k) f(x, a)|$  and  $f \mapsto \frac{1}{H} \int_{\mathcal{X}} |f(y)| d\bar{P}_h^k(y|x, a)$  are 2-Lipschitz with respect to  $\|\cdot\|_\infty$ . Consequently, with probability at least

$$1 - \delta K H N \left( \frac{\sigma^{2+d_2}}{H^2K}, \mathcal{X} \times \mathcal{A}, \rho \right) \left( \frac{2H}{L\sigma} \right)^{\mathcal{N}(\sigma, \mathcal{X}, \rho, x)},$$

for all  $\mathcal{L}(2L, 2H)$  and for all  $(x, a, h, k)$ , we have

$$\begin{aligned} \left| (\hat{P}_h^k - \bar{P}_h^k) f(x, a) \right| &\leq \frac{1}{H} \mathbb{P}_h^k |f(y)|(x, a) + \frac{(4H^2 + 10H)\Delta(k, \delta)}{\mathbf{C}_h^k(x, a)} \\ &\quad + \frac{2}{H} \mathbf{bias}_p(k, h) + \left( 2L_p L + \frac{4HC_2k}{\beta\sigma} + \frac{(4H^2 + 10H)\Delta(k, \delta)C_2k}{\beta^2\sigma} \right) \frac{\sigma^{2+d_2}}{H^2K} \\ &\quad + 6\sigma L_p L \left( 1 + \sqrt{\log^+(C_1 k / \beta)} \right) + \frac{2\beta H}{\mathbf{C}_h^k(x, a)} + 32L\sigma \end{aligned}$$

which concludes the proof.  $\square$

## D.4 Good event

**Lemma 9.** Let  $\mathcal{G} = \mathcal{G}_1 \cap \mathcal{G}_2 \cap \mathcal{G}_3 \cap \mathcal{G}_4$ , where

$$\begin{aligned} \mathcal{G}_1 &\stackrel{\text{def}}{=} \left\{ \forall(x, a, k, h), \left| \widehat{r}_h^k(x, a) - \bar{r}_h^k(x, a) \right| \leq \sqrt{\frac{2\Box_1^i(k, \delta/8)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_r(k, \delta/8)\sigma \right\} \\ \mathcal{G}_2 &\stackrel{\text{def}}{=} \left\{ \forall(x, a, k, h), \left| (\widehat{P}_h^k - \bar{P}_h^k)V_{k, h+1}^*(x, a) \right| \leq \sqrt{\frac{2H^2\Box_1^p(k, \delta/8)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} + \mathbf{b}_p(k, \delta/8)\sigma \right\} \\ \mathcal{G}_3 &\stackrel{\text{def}}{=} \left\{ \forall(x, a, k, h, f), \left| (\widehat{P}_h^k - \bar{P}_h^k)f(x, a) \right| \leq \sqrt{\frac{2H^2\Box_2^p(k, \delta/8)}{\mathbf{C}_h^k(x, a)}} + \frac{\beta H}{\mathbf{C}_h^k(x, a)} \right. \\ &\quad \left. + \theta_b^1(k, \delta/8)\sigma^{1+d_2/2} + \theta_b^2(k, \delta/8)\sigma \right\} \\ \mathcal{G}_4 &\stackrel{\text{def}}{=} \left\{ \forall(x, a, k, h, f), \left| (\widehat{P}_h^k - \bar{P}_h^k)f(x, a) \right| \leq \frac{1}{H}P_h^k |f(y)|(x, a) + \frac{14H^2C_2\Box_3(k, \delta/8) + 2\beta H}{\mathbf{C}_h^k(x, a)} \right. \\ &\quad \left. + \theta_b^3(k, \delta/8)\sigma^{1+d_2} + \theta_b^4(k, \delta/8)\sigma + \frac{2}{H} \mathbf{bias}_p(k, h) \right\} \end{aligned}$$

for  $(x, a, k, h) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$  and  $f \in \mathcal{L}(2L, 2H)$ , and where

$$\begin{aligned} \Box_1^i(k, \delta) &= \tilde{\mathcal{O}}(d_1), \quad \mathbf{b}_p(k, \delta) = \tilde{\mathcal{O}}(L + \sqrt{d_1}), \quad \Box_1^i(k, \delta) = \tilde{\mathcal{O}}(d_1), \quad \mathbf{b}_r(k, \delta) = \tilde{\mathcal{O}}(L + \sqrt{d_1}) \\ \Box_2^p(k, \delta) &= \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1d_2), \quad \theta_b^1(k, \delta) = \tilde{\mathcal{O}}(\sqrt{|\mathcal{C}'_\sigma|} + \sqrt{d_1d_2}), \quad \theta_b^2(k, \delta) = \tilde{\mathcal{O}}(L) \\ \Box_3(k, \delta) &= \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1d_2), \quad \theta_b^3(k, \delta) = \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1d_2 + L\sigma), \quad \theta_b^4(k, \delta) = \tilde{\mathcal{O}}(L) \end{aligned}$$

are defined in Lemmas 5, 6, 7 and 8, respectively. Then,

$$\mathbb{P}[\mathcal{G}] \geq 1 - \delta/2.$$

*Proof.* Immediate consequence of Lemmas 5, 6, 7 and 8. □

## E Upper bound on true value function

In this section, we show that the true value functions can be upper bounded by the value functions computed by **KeRNS** plus a bias term. This result will be used to upper bound the regret in the next section.

**Lemma 10** (upper bound on  $Q$  functions). *On  $\mathcal{G}$ , for all  $(x, a, k, h) \in \mathcal{X} \times \mathcal{A} \times [K] \times [H]$ , we have*

$$Q_h^k(x, a) + \sum_{h'=h}^H \mathbf{bias}(k, h) \geq Q_{k,h}^*(x, a)$$

where  $\mathbf{bias}(k, h) = \mathbf{bias}_r(k, h) + \mathbf{bias}_p(k, h)$  is the temporal bias of the algorithm at time  $(k, h)$  (see Definition 6).

*Proof.* We proceed by induction on  $h$ . For  $h = H + 1$ , both quantities are zero, so the inequality is trivially verified. Now, assume that it is true for  $h + 1$  and let's prove it for  $h$ .

From the induction hypothesis, we have  $V_{h+1}^k(x) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \geq V_{k,h+1}^*(x)$ . Indeed,

$$\max_a Q_{h+1}^k(x, a) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \geq \max_a Q_{k,h}^*(x, a) = V_{k,h+1}^*(x)$$

and, since  $V_{k,h+1}^*(x) \leq H - h$ , we have

$$V_{h+1}^k(x) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) = \min\left(H - h, \max_a Q_{h+1}^k(x, a)\right) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \geq V_{k,h+1}^*(x).$$

From the definition of the algorithm, we have

$$Q_h^k(x, a) = \min_{s \in [k-1]} \left[ \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho[(x, a), (x_h^s, a_h^s)] \right]$$

where  $\tilde{Q}_h^k(x, a) = \hat{r}_h^k(x, a) + \hat{P}_h^k V_{h+1}^k(x, a) + \mathbf{B}_h^k(x, a)$ . Hence,

$$\begin{aligned} & \tilde{Q}_h^k(x, a) - Q_{k,h}^*(x, a) \\ &= \underbrace{\hat{r}_h^k(x, a) - r_h^k(x, a) + r \mathbf{B}_h^k(x, a)}_{\text{(A)}} + \underbrace{\hat{P}_h^k V_{h+1}^k(x, a) - P_h^k V_{k,h+1}^*(x, a) + p \mathbf{B}_h^k(x, a)}_{\text{(B)}}. \end{aligned}$$

The term **(A)** is lower bounded as follows

$$\text{(A)} = \hat{r}_h^k(x, a) - \bar{r}_h^k(x, a) + r \mathbf{B}_h^k(x, a) + \bar{r}_h^k(x, a) - r_h^k(x, a) \geq -\mathbf{bias}_r(k, h)$$

by Corollary 2 and the fact that  $\hat{r}_h^k(x, a) - \bar{r}_h^k(x, a) + r \mathbf{B}_h^k(x, a) \geq 0$  on  $\mathcal{G}$ .

Similarly, for the term **(B)**, we have

$$\begin{aligned} \text{(B)} &= \hat{P}_h^k (V_{h+1}^k - V_{k,h+1}^*) (x, a) + \left( \hat{P}_h^k - \bar{P}_h^k \right) V_{k,h+1}^*(x, a) + \left( \bar{P}_h^k - P_h^k \right) V_{k,h+1}^*(x, a) + p \mathbf{B}_h^k(x, a) \\ &\geq \hat{P}_h^k (V_{h+1}^k - V_{k,h+1}^*) (x, a) - \mathbf{bias}_p(k, h) \end{aligned}$$

which gives us

$$\begin{aligned}
 & \tilde{Q}_h^k(x, a) - Q_{k,h}^*(x, a) \\
 & \geq \hat{P}_h^k(V_{h+1}^k - V_{k,h+1}^*)(x, a) - (\mathbf{bias}_r(k, h) + \mathbf{bias}_p(k, h)) \\
 & = \hat{P}_h^k(V_{h+1}^k - V_{k,h+1}^*)(x, a) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) - \sum_{h'=h}^H \mathbf{bias}(k, h) \\
 & \geq \hat{P}_h^k \left( V_{h+1}^k + \sum_{h'=h+1}^H \mathbf{bias}(k, h) - V_{k,h+1}^* \right) (x, a) - \sum_{h'=h}^H \mathbf{bias}(k, h) \\
 & \geq - \sum_{h'=h}^H \mathbf{bias}(k, h) \quad \text{by the induction hypothesis.}
 \end{aligned}$$

Consequently, for all  $(x, a)$  and all  $s \in [k-1]$ , we have

$$\begin{aligned}
 Q_{k,h}^*(x, a) - \sum_{h'=h}^H \mathbf{bias}(k, h) & \leq Q_{k,h}^*(x_h^s, a_h^s) + L\rho[(x, a), (x_h^s, a_h^s)] - \sum_{h'=h}^H \mathbf{bias}(k, h) \\
 & \leq \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho[(x, a), (x_h^s, a_h^s)]
 \end{aligned}$$

since  $Q_{k,h}^*$  is  $L$ -Lipschitz. Which implies the result

$$Q_{k,h}^*(x, a) - \sum_{h'=h}^H \mathbf{bias}(k, h) \leq \min_{s \in [k-1]} \left[ \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho[(x, a), (x_h^s, a_h^s)] \right] = Q_h^k(x, a).$$

□

**Corollary 3.** Let  $Q_{k,h}^+$  and  $V_{k,h}^+$  be defined as as

$$Q_{k,h}^+(x, a) \stackrel{\text{def}}{=} Q_h^k(x, a) + \sum_{h'=h}^H \mathbf{bias}(k, h'), \quad V_{k,h}^+(x) \stackrel{\text{def}}{=} \min \left( H, \max_a Q_{k,h}^+(x, a) \right)$$

Then,

$$\sup_{x \in \mathcal{X}} \left| V_h^k(x) - V_{k,h}^+(x) \right| \leq \sum_{h'=h}^H \mathbf{bias}(k, h')$$

and, by Lemma 10, we have  $V_{k,h}^+ \geq V_{k,h}^*$  on the event  $\mathcal{G}$ .

*Proof.* For any  $x \in \mathcal{X}$ ,

$$\begin{aligned}
 \left| V_h^k(x) - V_{k,h}^+(x) \right| & = \left| \min \left( H, \max_a Q_h^k(x, a) \right) - \min \left( H, \max_a Q_{k,h}^+(x, a) \right) \right| \\
 & \leq \left| \max_a Q_h^k(x, a) - \max_a Q_{k,h}^+(x, a) \right| \\
 & \leq \max_a \left| Q_h^k(x, a) - Q_{k,h}^+(x, a) \right| \\
 & = \sum_{h'=h}^H \mathbf{bias}(k, h').
 \end{aligned}$$

where we used the fact that, for any  $a, b, c \in \mathbb{R}$ , we have  $|\min(a, b) - \min(a, c)| \leq |b - c|$ .

□

## F Regret bounds

Using the results proved in the previous sections, we are now ready to prove our regret bounds. We first start by proving that the regret is bounded by sums involving  $\sqrt{1/\mathbf{C}_h^k(x, a)}$ ,  $1/\mathbf{C}_h^k(x, a)$  and bias terms. Then, we provide upper bounds for these sums, which result in the final regret bounds.

In Theorem 1, we prove two regret bounds,  $\mathcal{R}_1$  and  $\mathcal{R}_2$ . Here, we refer to  $\mathcal{R}_1$  as a UCRL-type regret bound and to  $\mathcal{R}_2$  as a UCBVI-type bound, due to the technique used to bound the difference between  $\widehat{P}_h^k$  and  $P_h^k$ . Making an analogy with finite MDPs, in UCRL (Jaksch et al., 2010), a term analogous to  $\|\widehat{P}_h - P_h^k\|_1$  is bounded (as in Lemma 7), whereas in UCBVI (Azar et al., 2017), the term  $|(\widehat{P}_h - P_h^k)V_{k,h+1}^*|$  is bounded (as in Lemma 5).

**Corollary 4.** Let  $\delta_h^k \stackrel{\text{def}}{=} V_h^k(x_h^k) - V_{k,h}^{\pi_k}(x_h^k)$ . Then, on  $\mathcal{G}$

$$\mathcal{R}(K) \leq \sum_{k=1}^K \delta_1^k + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h).$$

*Proof.* It follows directly from Lemma 10:

$$\mathcal{R}(K) = \sum_{k=1}^K V_{k,1}^*(x_1^k) - V_{k,h}^{\pi_k}(x_1^k) \leq \sum_{k=1}^K \left( V_1^k(x_1^k) + \sum_{h=1}^H \mathbf{bias}(k, h) - V_{k,h}^{\pi_k}(x_1^k) \right).$$

□

**Definition 8.** For any  $(k, h)$ , let  $(\tilde{x}_h^k, \tilde{a}_h^k)$  be defined as

$$(\tilde{x}_h^k, \tilde{a}_h^k) \stackrel{\text{def}}{=} \underset{(x_h^s, a_h^s): s < k}{\operatorname{argmin}} \rho[(x_h^k, a_h^k), (x_h^s, a_h^s)]$$

that is, the state-action pair in the history that is the closest to  $(x_h^k, a_h^k)$ .

### F.1 Regret bound in terms of the sum of exploration bonuses (UCRL-type)

**Lemma 11** (UCRL-type bound with sum of bonuses). On the event  $\mathcal{G}$ , the regret of *KeRNS* is bounded by

$$\begin{aligned} \mathcal{R}(K) \lesssim & \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H \sqrt{|\mathcal{C}'_\sigma|}}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I} \{ \rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} + H^2 |\mathcal{C}_\sigma| \\ & + \sum_{k=1}^K \sum_{h=1}^H \tilde{\xi}_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH\sigma \end{aligned}$$

where  $|\mathcal{C}_\sigma|$  is the  $\sigma$ -covering number of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $|\mathcal{C}'_\sigma|$  is the  $\sigma$ -covering number of  $(\mathcal{X}, \rho_{\mathcal{X}})$  and  $(\tilde{\xi}_{h+1}^k)_{k,h}$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .

*Proof.* **Regret decomposition** On  $\mathcal{G}$ , we upper bound  $\delta_h^k$  using the following decomposition:

$$\begin{aligned}
 \delta_h^k &= V_h^k(x_h^k) - V_{k,h}^{\pi_k}(x_h^k) \\
 &\leq Q_h^k(x_h^k, a_h^k) - Q_{k,h}^{\pi_k}(x_h^k, a_h^k) \\
 &\leq Q_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - Q_{k,h}^{\pi_k}(x_h^k, a_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)], \quad \text{since } Q_h^k \text{ is } L\text{-Lipschitz} \\
 &\leq \tilde{Q}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - Q_{k,h}^{\pi_k}(x_h^k, a_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)], \quad \text{since } Q_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \leq \tilde{Q}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \\
 &= \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k) + \hat{P}_h^k V_{h+1}^k(\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k V_{k,h+1}^{\pi_k}(x_h^k, a_h^k) + \mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)] \\
 &= \underbrace{\hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k)}_{\text{(A)}} + \underbrace{[\hat{P}_h^k - P_h^k] V_{k,h+1}^*(\tilde{x}_h^k, \tilde{a}_h^k)}_{\text{(B)}} + \underbrace{[\hat{P}_h^k - P_h^k] (V_{h+1}^k - V_{k,h+1}^*)}_{\text{(C)}}(\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\quad + \underbrace{P_h^k V_{h+1}^k(\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k V_{k,h+1}^{\pi_k}(x_h^k, a_h^k)}_{\text{(D)}} + \mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + 2L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)].
 \end{aligned}$$

Now, we bound each term (A)-(D) separately.

Term (A):

$$\begin{aligned}
 \text{(A)} &= \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k) \\
 &\leq \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + L_{\mathbf{r}}\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)] \\
 &= \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - \bar{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \bar{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + L_{\mathbf{r}}\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)] \\
 &\leq r \mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \mathbf{bias}_{\mathbf{r}}(k, h) + L_{\mathbf{r}}\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)]
 \end{aligned}$$

by the definition of  $\mathcal{G}$  and Corollary 2.

Term (B):

$$\begin{aligned}
 \text{(B)} &= [\hat{P}_h^k - P_h^k] V_{k,h+1}^*(\tilde{x}_h^k, \tilde{a}_h^k) = [\hat{P}_h^k - \bar{P}_h^k] V_{k,h+1}^*(\tilde{x}_h^k, \tilde{a}_h^k) + [\bar{P}_h^k - P_h^k] V_{k,h+1}^*(\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\leq {}^p \mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \mathbf{bias}_{\mathbf{p}}(k, h).
 \end{aligned}$$

Term (C): Using Corollary 2, we obtain

$$\begin{aligned}
 \text{(C)} &= [\hat{P}_h^k - P_h^k] (V_{h+1}^k - V_{k,h+1}^*) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\leq [\hat{P}_h^k - \bar{P}_h^k] (V_{h+1}^k - V_{k,h+1}^*) (\tilde{x}_h^k, \tilde{a}_h^k) + 2 \mathbf{bias}_{\mathbf{p}}(k, h) \\
 &\leq \sqrt{\frac{8H^2 \square_2^{\mathbf{p}}(k, \delta/8)}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{2\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + \theta_{\mathbf{b}}^1(k, \delta/8) \sigma^{1+d_2/2} + \theta_{\mathbf{b}}^2(k, \delta/8) \sigma + 2 \mathbf{bias}_{\mathbf{p}}(k, h) \\
 &\lesssim \sqrt{\frac{H^2 |\mathbf{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + 2 \mathbf{bias}_{\mathbf{p}}(k, h)
 \end{aligned}$$

by the definition of  $\mathcal{G}$ .

Term (D): From Assumption 2, for any  $L$ -Lipschitz function, the mapping  $(x, a) \mapsto P_h^k f(x, a)$  is  $L_{\mathbf{p}} L$ -Lipschitz. Consequently,

$$\begin{aligned}
 \text{(D)} &= P_h^k V_{h+1}^k(\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k V_{k,h+1}^{\pi_k}(x_h^k, a_h^k) \\
 &\leq P_h^k V_{h+1}^k(x_h^k, a_h^k) - P_h^k V_{k,h+1}^{\pi_k}(x_h^k, a_h^k) + L_{\mathbf{p}} L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &= P_h^k (V_{h+1}^k - V_{k,h+1}^{\pi_k})(x_h^k, a_h^k) + L_{\mathbf{p}} L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &= \delta_{h+1}^k + \xi_{h+1}^k + L_{\mathbf{p}} L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)]
 \end{aligned}$$

where

$$\xi_{h+1}^k \stackrel{\text{def}}{=} \mathbf{P}_h^k \left( V_{h+1}^k - V_{k,h+1}^{\pi_k} \right) (x_h^k, a_h^k) - \delta_{h+1}^k$$

is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .

Putting together the bounds for **(A)**-**(D)** and using the definition of the bonuses  $\mathbf{B}_h^k$ , we obtain

$$\delta_h^k \lesssim \delta_{h+1}^k + \xi_{h+1}^k + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] + \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \mathbf{bias}(k, h)$$

where the constant in front of  $\delta_{h+1}^k$  is *exact* (i.e., not omitted by  $\lesssim$ ).

Let  $E_h^k \stackrel{\text{def}}{=} \{\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$ . The inequality above implies

$$\mathbb{I}\{E_h^k\} \delta_h^k \lesssim \mathbb{I}\{E_h^k\} \delta_{h+1}^k + \mathbb{I}\{E_h^k\} \left( \xi_{h+1}^k + \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) + 3L\sigma + \mathbf{bias}(k, h). \quad (2)$$

Now, we bound  $\mathbb{I}\{E_h^k\} \delta_{h+1}^k$  in terms of  $\delta_{h+1}^k$ , which will be later used to bound  $\delta_h^k$  in terms of  $\delta_{h+1}^k$ . On  $\mathcal{G}$ , we have

$$\begin{aligned} \mathbb{I}\{E_h^k\} \delta_{h+1}^k &= \mathbb{I}\{E_h^k\} \left( V_{h+1}^k(x_{h+1}^k) - V_{k,h+1}^{\pi_k}(x_{h+1}^k) \right) \\ &= \mathbb{I}\{E_h^k\} \left( \underbrace{V_{h+1}^k(x_{h+1}^k) + \sum_{h'=h+1}^H \mathbf{bias}(k, h') - V_{k,h+1}^{\pi_k}(x_{h+1}^k)}_{\geq 0 \text{ by Lemma 10}} - \sum_{h'=h+1}^H \mathbf{bias}(k, h') \right) \\ &\leq V_{h+1}^k(x_{h+1}^k) + \sum_{h'=h+1}^H \mathbf{bias}(k, h') - V_{k,h+1}^{\pi_k}(x_{h+1}^k) - \mathbb{I}\{E_h^k\} \sum_{h'=h+1}^H \mathbf{bias}(k, h') \\ &= \delta_{h+1}^k + \sum_{h'=h+1}^H \mathbf{bias}(k, h') - \mathbb{I}\{E_h^k\} \sum_{h'=h+1}^H \mathbf{bias}(k, h') \\ &\leq \delta_{h+1}^k + \sum_{h'=h+1}^H \mathbf{bias}(k, h'). \end{aligned}$$

The inequality above, combined with (2) yields

$$\mathbb{I}\{E_h^k\} \delta_h^k \lesssim \delta_{h+1}^k + \sum_{h'=h}^H \mathbf{bias}(k, h') + \mathbb{I}\{E_h^k\} \left( \xi_{h+1}^k + \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) + 3L\sigma.$$

Let  $\overline{E}_h^k$  be the complement of  $E_h^k$ . Since  $\delta_h^k \leq H$ , we have

$$\begin{aligned} \delta_h^k &= \mathbb{I}\{\overline{E}_h^k\} \delta_h^k + \mathbb{I}\{E_h^k\} \delta_h^k \\ &\leq H \mathbb{I}\{\overline{E}_h^k\} + \mathbb{I}\{E_h^k\} \delta_h^k \\ &\lesssim H \mathbb{I}\{\overline{E}_h^k\} + \delta_{h+1}^k + \sum_{h'=h}^H \mathbf{bias}(k, h') + \mathbb{I}\{E_h^k\} \left( \xi_{h+1}^k + \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) + L\sigma. \end{aligned}$$



This yields

$$\begin{aligned} \delta_1^k &\lesssim \sum_{h=1}^H \mathbb{I}\{E_h^k\} \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \sum_{h=1}^H \mathbb{I}\{E_h^k\} \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\ &\quad + \sum_{h=1}^H \mathbb{I}\{E_h^k\} \xi_{h+1}^k + H \sum_{h=1}^H \mathbf{bias}(k, h) + H \sum_{h=1}^H \mathbb{I}\{\bar{E}_h^k\} + HL\sigma \end{aligned}$$

Using Corollary 4, we obtain

$$\begin{aligned} \mathcal{R}(K) &\leq \sum_{k=1}^K \delta_1^k + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) \\ &\lesssim \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \sqrt{\frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\ &\quad + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \xi_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + H \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{\bar{E}_h^k\} + KHL\sigma. \end{aligned}$$

For each  $h$ , the number of episodes  $k$  where the event  $\{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] > 2\sigma\}$  occurs is bounded by  $|\mathcal{C}_\sigma|$ . Hence, we can bound the sum

$$H \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{\bar{E}_h^k\} = H \sum_{h=1}^H \sum_{k=1}^K \mathbb{I}\{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] > 2\sigma\} \leq H^2 |\mathcal{C}_\sigma|.$$

We conclude the proof by recalling the definition  $E_h^k \stackrel{\text{def}}{=} \{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$  and using the fact that  $\tilde{\xi}_{h+1}^k \stackrel{\text{def}}{=} \mathbb{I}\{E_h^k\} \xi_{h+1}^k$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .  $\square$

## F.2 Regret bound in terms of the sum of exploration bonuses (UCBVI-type)

**Lemma 12** (UCBVI-type bound with sum of bonuses). *In the event  $\mathcal{G}$ , the regret of  $\text{KeRMS}$  is bounded by*

$$\begin{aligned} \mathcal{R}(K) &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I}\{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\} + H^2 |\mathcal{C}_\sigma| \\ &\quad + \sum_{k=1}^K \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^h \tilde{\xi}_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH\sigma \end{aligned}$$

where  $|\mathcal{C}_\sigma|$  is the  $\sigma$ -covering number of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $|\mathcal{C}'_\sigma|$  is the  $\sigma$ -covering number of  $(\mathcal{X}, \rho_{\mathcal{X}})$  and  $(\tilde{\xi}_{h+1}^k)_{k,h}$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .

*Proof.* The proof follows the one of Proposition 5 of Domingues et al. (2020). The key difference is that we need to handle the temporal bias. In particular,  $V_h^k$  is not an upper bound on  $V_{k,h+1}^*$  due to the temporal bias, which makes our proof slightly more technical by introducing  $V_{k,h}^+$  (see Cor. 3) when applying the Bernstein-type concentration of Lemma 8.

**Regret decomposition** We use the same regret decomposition as in the proof of Lemma 11. The terms (A), (B) and (D) are bounded in the same way, but we handle the term (C) differently.

Term (C): To bound this term, we use corollaries 2 and 3:

$$\begin{aligned}
 \text{(C)} &= \left[ \widehat{P}_h^k - P_h^k \right] \left( V_{h+1}^k - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &= \left[ \widehat{P}_h^k - P_h^k \right] \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \widehat{P}_h^k - P_h^k \right] \left( V_{h+1}^k - V_{k,h+1}^+ \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\leq \left[ \widehat{P}_h^k - P_h^k \right] \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h'), \quad \text{by Cor. 3} \\
 &= \left[ \widehat{P}_h^k - \overline{P}_h^k \right] \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \overline{P}_h^k - P_h^k \right] \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h') \\
 &\leq \left[ \widehat{P}_h^k - \overline{P}_h^k \right] \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + 2 \mathbf{bias}_p(k, h) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h'), \quad \text{by Cor. 2} \\
 &\leq \frac{1}{H} P_h^k \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \frac{14H^2 C_2 \square_3(k, \delta/8) + 2\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\
 &\quad + \theta_b^3(k, \delta/8) \sigma^{1+d_2} + \theta_b^4(k, \delta/8) \sigma + \frac{2}{H} \mathbf{bias}_p(k, h) + 2 \mathbf{bias}_p(k, h) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h') \\
 &\leq \frac{1}{H} P_h^k \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (x_h^k, a_h^k) + \frac{14H^2 C_2 \square_3(k, \delta/8) + 2\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + \frac{2L_p L}{H} \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &\quad + \theta_b^3(k, \delta/8) \sigma^{1+d_2} + \theta_b^4(k, \delta/8) \sigma + \frac{2}{H} \mathbf{bias}_p(k, h) + 2 \mathbf{bias}_p(k, h) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h'),
 \end{aligned}$$

where we also used the definition of  $\mathcal{G}$  and the fact that the function  $(x, a) \mapsto P_h^k \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (x, a)$  is  $2L_p L$ -Lipschitz, from Assumption 2. Now, since

$$\square_3(k, \delta) = \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1 d_2), \quad \theta_b^3(k, \delta) = \tilde{\mathcal{O}}(|\mathcal{C}'_\sigma| + d_1 d_2 + L\sigma), \quad \theta_b^4(k, \delta) = \tilde{\mathcal{O}}(L)$$

and  $|\mathcal{C}'_\sigma| = \mathcal{O}(1/\sigma^{d_2})$ , we have

$$\begin{aligned}
 \text{(C)} &\lesssim \frac{1}{H} P_h^k \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (x_h^k, a_h^k) + \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma \\
 &\quad + \left( 2 + \frac{2}{H} \right) \mathbf{bias}_p(k, h) + 2 \sum_{h'=h}^H \mathbf{bias}(k, h') + \frac{2L_p L}{H} \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)].
 \end{aligned}$$

Using again Corollary 3, we have

$$\frac{1}{H} P_h^k \left( V_{k,h+1}^+ - V_{k,h+1}^* \right) (x_h^k, a_h^k) \leq \frac{1}{H} P_h^k \left( V_{h+1}^k - V_{k,h+1}^* \right) (x_h^k, a_h^k) + \frac{1}{H} \sum_{h'=h}^H \mathbf{bias}(k, h')$$

which gives us, since  $V_{k,h+1}^{\pi_k} \leq V_{k,h+1}^*$ ,

$$\begin{aligned}
 \text{(C)} &\lesssim \frac{1}{H} P_h^k \left( V_{h+1}^k - V_{k,h+1}^* \right) (x_h^k, a_h^k) + \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \sum_{h'=h}^H \mathbf{bias}(k, h') + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &\lesssim \frac{1}{H} P_h^k \left( V_{h+1}^k - V_{k,h+1}^{\pi_k} \right) (x_h^k, a_h^k) + \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \sum_{h'=h}^H \mathbf{bias}(k, h') + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)]
 \end{aligned}$$

where we omit constants. Notice, however, that there are no constants omitted in the term  $\frac{1}{H} P_h^k \left( V_{h+1}^k - V_{k,h+1}^{\pi_k} \right) (x_h^k, a_h^k)$ .

Putting together the bounds for (A)-(D), we obtain

$$\delta_h^k \lesssim \left(1 + \frac{1}{H}\right) (\delta_{h+1}^k + \xi_{h+1}^k) + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] + 2\mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \sum_{h'=h}^H \mathbf{bias}(k, h')$$

where the constant in front of  $\delta_{h+1}^k$  is *exact* (i.e., not omitted by  $\lesssim$ ).

Let  $E_h^k \stackrel{\text{def}}{=} \{\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$ . Using the definition of the bonus

$$\mathbf{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \lesssim \frac{H}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma,$$

and the same argument as in the proof of Lemma 11, we obtain

$$\begin{aligned} \mathcal{R}(K) &\leq \sum_{k=1}^K \delta_1^k + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) \\ &\lesssim \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \sqrt{\frac{H^2}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \frac{H^2 |\mathcal{C}'_\sigma|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\ &\quad + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \xi_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + H^2 |\mathcal{C}_\sigma| + KH L\sigma. \end{aligned}$$

As in Lemma 11, we conclude the proof by recalling the definition  $E_h^k \stackrel{\text{def}}{=} \{\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$  and using the fact that  $\tilde{\xi}_{h+1}^k \stackrel{\text{def}}{=} \mathbb{I}\{E_h^k\} \xi_{h+1}^k$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .  $\square$

### F.3 Bounding the sum of bonuses and bias

**Lemma 13.** *Let  $(\mu_i)_{i \geq 1}$  be a sequence of non-negative numbers. Then,*

$$\sum_{k=1}^K \sum_{i=1 \vee (k-W)}^{k-1} \mu_i \leq 2W \sum_{i=1}^K \mu_i.$$

*Proof.* We have

$$\begin{aligned} \sum_{k=1}^K \sum_{i=1 \vee (k-W)}^{k-1} \mu_i &= \sum_{k=1}^W \sum_{i=1}^{k-1} \mu_i + \sum_{k=W+1}^K \sum_{i=k-W}^{k-1} \mu_i \leq W \sum_{i=1}^K \mu_i + \sum_{i=1}^{K-1} \sum_{k=i+1}^{i+W} \mu_i \\ &\leq W \sum_{i=1}^K \mu_i + W \sum_{i=1}^K \mu_i = 2W \sum_{i=1}^K \mu_i. \end{aligned}$$

$\square$

**Corollary 5** (bound on the temporal bias). *Let  $\Delta^r$  and  $\Delta^p$  be the variation of the MDP over  $KH$  time steps,*

$$\Delta^r \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} |r_h^i(x, a) - r_h^{i+1}(x, a)|, \quad \Delta^p \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} \mathbb{W}_1 (P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)).$$

*Then,*

$$\sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(h, k) \leq 2W (\Delta^r + L\Delta^p) + \frac{2C_3(H+1)KH}{\beta} \frac{\eta^W}{1-\eta}.$$

*Proof.* From Lemma 13, we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}_p(h, k) &= \sum_{k=1}^K \sum_{h=1}^H L \sum_{i=1 \vee (k-W)}^{k-1} \sup_{x, a} \mathbb{W}_1(P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)) + \sum_{k=1}^K \sum_{h=1}^H \frac{2C_3 H}{\beta} \frac{\eta^W}{1-\eta} \\ &\leq 2WL \sum_{h=1}^H \sum_{i=1}^K \sup_{x, a} \mathbb{W}_1(P_h^i(\cdot|x, a), P_h^{i+1}(\cdot|x, a)) + \frac{2C_3 KH^2}{\beta} \frac{\eta^W}{1-\eta}. \end{aligned}$$

The sum  $\sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}_r(h, k)$  is bounded in the same way, which concludes the proof.  $\square$

**Lemma 14** (bounding sum on sliding window). *Let  $\{a_n\}_{n \geq 1}$  be a sequence of real numbers such that  $0 \leq a \leq c$  for some constant  $c > 0$ . Let  $A_t = \sum_{n=1 \vee (t-W)}^{t-1} a_n$ . Then, for any  $p, b > 0$ ,*

$$\sum_{t=1}^T \frac{a_t}{(1 + bA_t)^p} \leq \sum_{n=1}^{\lceil T/W \rceil} \left( c + \int_0^{A_{nW+1-c}} \frac{1}{(1 + bz)^p} dz \right).$$

*Proof.* We have

$$\sum_{t=1}^T \frac{a_t}{(1 + bA_t)^p} = \underbrace{\sum_{t=1}^W \frac{a_t}{(1 + bA_t)^p}}_{\textcircled{1}} + \underbrace{\sum_{t=W+1}^T \frac{a_t}{(1 + bA_t)^p}}_{\textcircled{2}}.$$

By Lemma 9 of Domingues et al. (2020), we have

$$\textcircled{1} = \sum_{t=1}^W \frac{a_t}{(1 + bA_t)^p} \leq c + \int_0^{A_{W+1-c}} \frac{1}{(1 + bz)^p} dz.$$

Now, we handle  $\textcircled{2}$ :

$$\begin{aligned} \textcircled{2} &\leq \sum_{n=1}^{\lceil T/W \rceil - 1} \sum_{t=nW+1}^{(n+1)W} \frac{a_t}{(1 + bA_t)^p} = \sum_{n=1}^{\lceil T/W \rceil - 1} \sum_{l=1}^W \frac{a_{l+nW}}{(1 + bA_{l+nW})^p} \\ &\leq \sum_{n=1}^{\lceil T/W \rceil - 1} \left( c + \int_0^{A_{(n+1)W+1-c}} \frac{1}{(1 + bz)^p} dz \right). \end{aligned}$$

$\square$

**Definition 9.** Consider a  $\sigma$ -covering of  $(\mathcal{X} \times \mathcal{A}, \rho)$ ,  $\mathcal{C}_\sigma = \{(x_j, a_j) \in \mathcal{X} \times \mathcal{A}, j = 1, \dots, |\mathcal{C}_\sigma|\}$ . We define a partition  $\{B_j\}_{j \in [|\mathcal{C}_\sigma|]}$  such that

$$B_j = \left\{ (x, a) \in \mathcal{X} \times \mathcal{A} : (x_j, a_j) = \underset{(x_i, a_i) \in \mathcal{C}_\sigma}{\operatorname{argmin}} \rho[(x, a), (x_i, a_i)] \right\}$$

with ties broken arbitrarily.

**Lemma 15.** Let  $U_\eta \stackrel{\text{def}}{=} \lceil 1/\log(1/\eta) \rceil$  and

$$\mathbf{N}_h^k(B_j, U_\eta) \stackrel{\text{def}}{=} \sum_{s=1 \vee (k-U_\eta)}^{k-1} \mathbb{I}\{(x_h^s, a_h^s) \in B_j\}.$$

If  $U_\eta \leq W$ ,  $\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma$  and  $(x_h^k, a_h^k) \in B_j$  then

$$\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \geq \beta + G(4)e^{-1} \mathbf{N}_h^k(B_j, U_\eta).$$

*Proof.* This result is based on the proof of Proposition 6 of Domingues et al. (2020), which we generalize to the case where the kernel is time-dependent. From Assumption 4, we have  $\bar{\Gamma}_{(\sigma, \eta, W)}(k-1-s, z) \geq G(z)\eta^{k-1-s}$  for all  $s \geq k-W$ . Consequently, if  $U_\eta \leq W$ , for all  $s \geq k-U_\eta$ :

$$\bar{\Gamma}_{(\sigma, \eta, W)}(k-1-s, z) \geq G(z)\eta^{k-1-s} \geq G(z)\eta^{U_\eta} \geq G(z)\exp(-1). \quad (3)$$

Also, if  $\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma$ ,  $(x_h^k, a_h^k) \in B_j$ , and  $(x_h^s, a_h^s) \in B_j$ , we have

$$\begin{aligned} \rho[(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^s, a_h^s)] &\leq \rho[(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)] + \rho[(x_h^k, a_h^k), (x_h^s, a_h^s)] \\ &\leq 2\sigma + \rho[(x_h^k, a_h^k), (x_j, a_j)] + \rho[(x_j, a_j), (x_h^s, a_h^s)] \leq 4\sigma. \end{aligned} \quad (4)$$

By Assumption 4, the function  $z \mapsto \bar{\Gamma}_{(\eta, W)}(t, z)$  is non-increasing. Together with (3) and (4), this yields:

$$\begin{aligned} \mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) &= \beta + \sum_{s=1}^{k-1} \bar{\Gamma}_{(\sigma, \eta, W)}\left(k-1-s, \frac{\rho[(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^s, a_h^s)]}{\sigma}\right) \\ &\geq \beta + \sum_{s=1}^{k-1} \bar{\Gamma}_{(\sigma, \eta, W)}\left(k-1-s, \frac{\rho[(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^s, a_h^s)]}{\sigma}\right) \mathbb{I}\{(x_h^s, a_h^s) \in B_j\} \\ &\geq \beta + \sum_{s=1 \vee (k-U_\eta)}^{k-1} \bar{\Gamma}_{(\sigma, \eta, W)}(k-1-s, 4) \mathbb{I}\{(x_h^s, a_h^s) \in B_j\} \\ &\geq \beta + G(4)e^{-1} \sum_{s=k-U_\eta}^{k-1} \mathbb{I}\{(x_h^s, a_h^s) \in B_j\}, \end{aligned}$$

which concludes the proof.  $\square$

**Lemma 16.** Let  $U_\eta = \lceil 1/\log(1/\eta) \rceil$ . If  $U_\eta \leq W$ , we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} \mathbb{I}\{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\} &\lesssim H \left\lceil \frac{K}{U_\eta} \right\rceil \left( |\mathcal{C}_\sigma| + \sqrt{|\mathcal{C}_\sigma| U_\eta} \right) \\ \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \mathbb{I}\{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\} &\lesssim H |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil. \end{aligned}$$

*Proof.* The proof relies on Lemmas 14 and 15.

Here, we define the constant  $c$  as  $c = G(4)\beta^{-1}e^{-1} > 0$ , since  $G(4) > 0$  by Assumption 4.

Bounding the sum  $\sum_k 1/\sqrt{\mathbf{C}_h^k}$

$$\begin{aligned}
 & \sum_{k=1}^K \frac{1}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} \\
 &= \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{k=1}^K \frac{1}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} \mathbb{I} \{ (x_h^k, a_h^k) \in B_j \} \\
 &\leq \beta^{-1/2} \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{k=1}^K \frac{\mathbb{I} \{ (x_h^k, a_h^k) \in B_j \}}{\sqrt{1 + c\mathbf{N}_h^k(B_j, U_\eta)}}, \quad \text{by Lemma 15} \\
 &\leq \beta^{-1/2} \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{n=1}^{\lceil K/U_\eta \rceil} \left( 1 + \int_0^{\mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta)} \frac{1}{\sqrt{1 + cz}} dz \right), \quad \text{by Lemma 14} \\
 &= \beta^{-1/2} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{2\beta^{-1/2}}{c} \sum_{n=1}^{\lceil K/U_\eta \rceil} \sum_{j=1}^{|\mathcal{C}_\sigma|} \sqrt{1 + c\mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta)} \\
 &\leq \beta^{-1/2} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{2\beta^{-1/2}}{c} \sum_{n=1}^{\lceil K/U_\eta \rceil} \sqrt{|\mathcal{C}_\sigma|} \sqrt{|\mathcal{C}_\sigma| + c \sum_{j=1}^{|\mathcal{C}_\sigma|} \mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta)}, \quad \text{by Cauchy-Schwarz inequality} \\
 &\leq \beta^{-1/2} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{2\beta^{-1/2}}{c} \sum_{n=1}^{\lceil K/U_\eta \rceil} \sqrt{|\mathcal{C}_\sigma|} \sqrt{|\mathcal{C}_\sigma| + cU_\eta} \lesssim |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \sqrt{|\mathcal{C}_\sigma| U_\eta} \left\lceil \frac{K}{U_\eta} \right\rceil.
 \end{aligned}$$

Bounding the sum  $\sum_k 1/\mathbf{C}_h^k$

$$\begin{aligned}
 & \sum_{k=1}^K \frac{1}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} \\
 &= \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{k=1}^K \frac{1}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} \mathbb{I} \{ (x_h^k, a_h^k) \in B_j \} \\
 &\leq \beta^{-1} \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{k=1}^K \frac{\mathbb{I} \{ (x_h^k, a_h^k) \in B_j \}}{1 + c\mathbf{N}_h^k(B_j, U_\eta)}, \quad \text{by Lemma 15} \\
 &\leq \beta^{-1} \sum_{j=1}^{|\mathcal{C}_\sigma|} \sum_{n=1}^{\lceil K/U_\eta \rceil} \left( 1 + \int_0^{\mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta)} \frac{1}{1 + cz} dz \right), \quad \text{by Lemma 14} \\
 &= \beta^{-1} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{\beta^{-1} |\mathcal{C}_\sigma|}{c} \sum_{n=1}^{\lceil K/U_\eta \rceil} \sum_{j=1}^{|\mathcal{C}_\sigma|} \frac{1}{|\mathcal{C}_\sigma|} \log \left( 1 + c\mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta) \right) \\
 &\leq \beta^{-1} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{\beta^{-1} |\mathcal{C}_\sigma|}{c} \sum_{n=1}^{\lceil K/U_\eta \rceil} \log \left( 1 + \frac{c}{|\mathcal{C}_\sigma|} \sum_{j=1}^{|\mathcal{C}_\sigma|} \mathbf{N}_h^{nU_\eta+1}(B_j, U_\eta) \right), \quad \text{by Jensen's inequality} \\
 &\leq \beta^{-1} |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + \frac{\beta^{-1} |\mathcal{C}_\sigma|}{c} \left\lceil \frac{K}{U_\eta} \right\rceil \log \left( 1 + c \frac{U_\eta}{|\mathcal{C}_\sigma|} \right) \lesssim |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil.
 \end{aligned}$$

□

## F.4 Final regret bounds

### F.4.1 UCRL-type regret bounds

**Theorem 3** (UCRL-type regret bound). *If  $U_\eta = \lceil 1/\log(1/\eta) \rceil \leq W$ , the regret of  $\text{KeRNS}$  is bounded by*

$$\begin{aligned} \mathcal{R}(K) &\lesssim H^2 \left\lceil \frac{K}{U_\eta} \right\rceil \sqrt{|\mathcal{C}'_\sigma|} \left( |\mathcal{C}_\sigma| + \sqrt{|\mathcal{C}_\sigma| U_\eta} \right) + H^2 |\mathcal{C}_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + H^{3/2} \sqrt{K} \\ &\quad + W (\Delta^r + L\Delta^p) H + \frac{\eta^W}{1-\eta} K H^3 \\ &\quad + H^2 |\mathcal{C}_\sigma| + LKH\sigma \end{aligned}$$

with probability at least  $1 - \delta$ , where  $|\mathcal{C}_\sigma|$  and  $|\mathcal{C}'_\sigma|$  are the  $\sigma$ -covering numbers of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and  $(\mathcal{X}, \rho_{\mathcal{X}})$ , respectively, and

$$\Delta^r \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} |r_h^i(x,a) - r_h^{i+1}(x,a)|, \quad \Delta^p \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} \mathbb{W}_1(P_h^i(\cdot|x,a), P_h^{i+1}(\cdot|x,a))$$

represent the variation of the rewards and transitions, respectively.

*Proof.* We apply Lemma 11, Lemma 16 and Corollary 5 and the fact that  $\mathbb{P}[\mathcal{G}] \geq 1 - \delta/2$  by Lemma 9. To conclude, notice that  $\sum_{k=1}^K \sum_{h=1}^H \tilde{\xi}_{h+1}^k \lesssim H\sqrt{KH}$  with probability at least  $1 - \delta/2$  by Hoeffding-Azuma's inequality.  $\square$

**Corollary 6.** *Let  $d_1$  and  $d_2$  be the covering dimensions of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and  $(\mathcal{X}, \rho_{\mathcal{X}})$ , respectively. Let  $\alpha = \frac{1}{d_1+d_2+3}$ ,  $\Delta = \Delta^r + L\Delta^p$  and*

$$\sigma = K^{-\alpha}, \quad \log\left(\frac{1}{\eta}\right) = \left(\frac{\Delta}{K^{1+\alpha(d_1+d_2)/2}}\right)^{2/3}, \quad W = \frac{\log(K/(1-\eta))}{\log(1/\eta)}$$

Since  $W \geq U_\eta = \lceil 1/\log(1/\eta) \rceil$ , we have, with probability at least  $1 - \delta$ ,

$$\begin{aligned} \mathcal{R}(K) &\lesssim H^2 \Delta^{\frac{2}{3}} K^{\frac{d_1+d_2/6+1}{d_1+d_2+3}} + H^2 \Delta^{\frac{1}{3}} K^{\frac{d_1+d_2+2}{d_1+d_2+3}} \\ &\quad + H^2 \Delta^{\frac{2}{3}} K^{\frac{d_1+d_2/6+1}{d_1+d_2+3}} + H^{\frac{3}{2}} \sqrt{K} + H \Delta^{\frac{1}{3}} K^{\frac{d_1+d_2+2}{d_1+d_2+3}} \log\left(\frac{HK}{1-\eta}\right) + H^3 \\ &\quad + H^2 K^{\frac{d_1}{d_1+d_2+3}} + LHK K^{\frac{d_1+d_2+2}{d_1+d_2+3}} \end{aligned}$$

that is,

$$\mathcal{R}(K) = \tilde{\mathcal{O}}\left(\Delta^{\frac{1}{3}} K^{\frac{d_1+d_2+2}{d_1+d_2+3}} \left(H^2 + H \log\left(\frac{HK}{1-\eta}\right)\right)\right).$$

Furthermore, if  $\limsup_{K \rightarrow \infty} \Delta/K = 0$ , we have

$$\eta = \exp\left(-\Delta^{2/3}/K^{\frac{d_1+d_2+2}{d_1+d_2+3}}\right) \underset{K \rightarrow \infty}{\sim} 1 - \Delta^{2/3}/K^{\frac{d_1+d_2+2}{d_1+d_2+3}}$$

which implies

$$\mathcal{R}(K) = \tilde{\mathcal{O}}\left(H^2 \Delta^{\frac{1}{3}} K^{\frac{d_1+d_2+2}{d_1+d_2+3}}\right).$$

*Proof.* Immediate consequence of Theorem 3 and the fact that  $|\mathcal{C}_\sigma| = \mathcal{O}(\sigma^{-d_1})$  and  $|\mathcal{C}'_\sigma| = \mathcal{O}(\sigma^{-d_2})$ .  $\square$

**Corollary 7** (UCRL-type regret bound in discrete case). *If  $\mathcal{X} \times \mathcal{A}$  is finite, we can take  $\sigma = 0$  and  $|\mathcal{C}_\sigma| = XA$ ,  $|\mathcal{C}'_\sigma| = X$ , where  $X = |\mathcal{X}|$  and  $A = |\mathcal{A}|$ . In this case, Theorem 3 and Corollary 6 give us*

$$\mathcal{R}(K) = \tilde{\mathcal{O}}\left(H^2 X \sqrt{A} \Delta^{\frac{1}{3}} K^{\frac{2}{3}}\right).$$

#### F.4.2 UCBVI-type regret bounds

**Theorem 4** (UCBVI-type regret bound). *If  $U_\eta = \lceil 1/\log(1/\eta) \rceil \leq W$ , the regret of **KeRNS** is bounded by*

$$\begin{aligned} \mathcal{R}(K) &\lesssim H^2 \left\lceil \frac{K}{U_\eta} \right\rceil \left( |\mathcal{C}_\sigma| + \sqrt{|\mathcal{C}_\sigma| U_\eta} \right) + H^3 |\mathcal{C}_\sigma| |\mathcal{C}'_\sigma| \left\lceil \frac{K}{U_\eta} \right\rceil + H^{3/2} \sqrt{K} \\ &\quad + W (\Delta^r + L \Delta^p) H + \frac{\eta^W}{1-\eta} K H^3 \\ &\quad + H |\mathcal{C}_\sigma| + L K H \sigma \end{aligned}$$

with probability at least  $1 - \delta$ , where  $|\mathcal{C}_\sigma|$  and  $|\mathcal{C}'_\sigma|$  are the  $\sigma$ -covering numbers of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and  $(\mathcal{X}, \rho_{\mathcal{X}})$ , respectively, and

$$\Delta^r \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} |r_h^i(x,a) - r_h^{i+1}(x,a)|, \quad \Delta^p \stackrel{\text{def}}{=} \sum_{i=1}^K \sum_{h=1}^H \sup_{x,a} \mathbb{W}_1(P_h^i(\cdot|x,a), P_h^{i+1}(\cdot|x,a))$$

represent the variation of the rewards and transitions, respectively.

*Proof.* We apply Lemma 12, Lemma 16 and Corollary 5 and the fact that  $\mathbb{P}[\mathcal{G}] \geq 1 - \delta/2$  by Lemma 9. To conclude, notice that  $\sum_{k=1}^K \sum_{h=1}^H \left(1 + \frac{1}{H}\right)^h \tilde{\xi}_{h+1}^k \lesssim H \sqrt{KH}$  with probability at least  $1 - \delta/2$  by Hoeffding-Azuma's inequality.  $\square$



**Corollary 8** (UCBVI-type regret bound in discrete case). *If  $\mathcal{X} \times \mathcal{A}$  is finite, we can take  $\sigma = 0$  and  $|\mathcal{C}_\sigma| = XA$ ,  $|\mathcal{C}'_\sigma| = X$ , where  $X = |\mathcal{X}|$  and  $A = |\mathcal{A}|$ . In this case, Theorem 4 gives us*

$$\begin{aligned} \mathcal{R}(K) &\lesssim H^2 \left\lceil \frac{K}{U_\eta} \right\rceil \left( XA + \sqrt{XAU_\eta} \right) + H^3 X^2 A \left\lceil \frac{K}{U_\eta} \right\rceil + H^{3/2} \sqrt{K} \\ &\quad + W (\Delta^r + L\Delta^p) H + \frac{\eta^W}{1-\eta} KH^3 + H^2 XA \end{aligned}$$

Let  $\Delta = \Delta^r + L\Delta^p$ . By choosing

$$\log\left(\frac{1}{\eta}\right) = \left(\frac{\Delta}{K}\right)^{2/3}, \quad W = \frac{\log(K/(1-\eta))}{\log(1/\eta)}$$

we obtain

$$\begin{aligned} \mathcal{R}(K) &\lesssim H^2 XA \Delta^{\frac{2}{3}} K^{\frac{1}{3}} + H^2 \sqrt{XA} \Delta^{\frac{1}{3}} K^{\frac{2}{3}} + H^3 X^2 A \Delta^{\frac{2}{3}} K^{\frac{1}{3}} \\ &\quad + H^{3/2} \sqrt{K} + \log\left(\frac{K}{1-\eta}\right) H \Delta^{\frac{1}{3}} K^{\frac{2}{3}} + H^3 + HXA. \end{aligned}$$

since  $W \geq U_\eta = \lceil 1/\log(1/\eta) \rceil$ . Furthermore, if  $\limsup_{K \rightarrow \infty} \Delta/K = 0$ , we have

$$\eta = \exp\left(-\Delta^{\frac{2}{3}}/K^{\frac{2}{3}}\right) \underset{K \rightarrow \infty}{\sim} 1 - \Delta^{\frac{2}{3}}/K^{\frac{2}{3}}$$

which implies

$$\mathcal{R}(K) = \tilde{\mathcal{O}}\left(H^2 \sqrt{XA} \Delta^{\frac{1}{3}} K^{\frac{2}{3}} + H^3 X^2 A \Delta^{\frac{2}{3}} K^{\frac{1}{3}}\right).$$

**Corollary 9.** *Let  $d_1$  and  $d_2$  be the covering dimensions of  $(\mathcal{X} \times \mathcal{A}, \rho)$  and  $(\mathcal{X}, \rho_X)$ , respectively. Let  $\alpha = \frac{1}{d_1+d_2+2}$ ,  $\Delta = \Delta^r + L\Delta^p$  and*

$$\sigma = K^\alpha, \quad \log\left(\frac{1}{\eta}\right) = \left(\frac{\Delta}{HK^{1+\alpha(d_1+d_2)}}\right)^{1/2}, \quad W = \left\lceil \frac{\log(K/(1-\eta))}{\log(1/\eta)} \right\rceil$$

Since  $W \geq U_\eta = \lceil 1/\log(1/\eta) \rceil$ , we have, with probability at least  $1 - \delta$ ,

$$\mathcal{R}(K) \lesssim H^2 \left(1 + \log\left(\frac{K}{1-\eta}\right)\right) \Delta^{\frac{1}{2}} K^{\frac{d_1+d_2+1}{d_1+d_2+2}} + H^{\frac{3}{2}} \Delta^{\frac{1}{4}} K^{\frac{3}{4}} + LHK^{\frac{d_1+d_2+1}{d_1+d_2+2}} + H^2.$$

Furthermore, if  $\limsup_{K \rightarrow \infty} \Delta/K = 0$ , we have

$$\eta = \exp\left(-\Delta^{1/2}/K^{\frac{d_1+d_2+1}{d_1+d_2+2}}\right) \underset{K \rightarrow \infty}{\sim} 1 - \Delta^{1/2}/K^{\frac{d_1+d_2+1}{d_1+d_2+2}}$$

which implies

$$\mathcal{R}(K) = \tilde{\mathcal{O}}\left(H^2 \Delta^{\frac{1}{2}} K^{\frac{d_1+d_2+1}{d_1+d_2+2}} + H^{\frac{3}{2}} \Delta^{\frac{1}{4}} K^{\frac{3}{4}}\right).$$

*Proof.* Immediate consequence of Theorem 4 and the fact that  $|\mathcal{C}_\sigma| = \mathcal{O}(\sigma^{-d_1})$  and  $|\mathcal{C}'_\sigma| = \mathcal{O}(\sigma^{-d_2})$ .  $\square$

## G RS-KeRNS: An efficient version of KeRNS using representative states

RS-KeRNS is described in Algorithm 4, which uses a backward induction on representative states (Algorithm 5) and updates the model online (algorithms 6 and 7). In this section, we introduce the main definitions used by RS-KeRNS, and we analyze its runtime and regret.

### G.1 Definitions

In each episode  $k$  and for each  $h$ , RS-KeRNS keeps and updates sets of representative states  $\bar{\mathcal{X}}_h^k$ , actions  $\bar{\mathcal{A}}_h^k$ , and next-states  $\bar{\mathcal{Y}}_h^k$ , with cardinalities  $\bar{X}_h^k$ ,  $\bar{A}_h^k$  and  $\bar{Y}_h^k$ , respectively. These sets are built using the data observed up to episode  $k - 1$ . We define the following projections:

$$\zeta_h^{k+1}(x, a) \stackrel{\text{def}}{=} \underset{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k}{\text{argmin}} \rho[(x, a), (\bar{x}, \bar{a})], \quad \bar{\zeta}_h^{k+1}(y) \stackrel{\text{def}}{=} \underset{\bar{y} \in \bar{\mathcal{Y}}_h^k}{\text{argmin}} \rho_{\mathcal{X}}(y, \bar{y}).$$

where we also assume to have access to the metric  $\rho_{\mathcal{X}}$ . The definitions below introduce the kernel function and the estimated MDP used by RS-KeRNS.

**Definition 10** (kernel function for RS-KeRNS). *Let  $(\alpha_i)_{i \geq 1}$  be a sequence of numbers in  $[0, 1]$ . RS-KeRNS uses a kernel of the form  $\Gamma(t, u, v) = \chi(t)\phi(u, v)$ , where*

$$\chi(t) \stackrel{\text{def}}{=} \prod_{i=1}^t \alpha_i, \quad \phi(u, v) \stackrel{\text{def}}{=} \exp\left(-\rho[u, v]^2 / (2\sigma^2)\right),$$

and, by convention,  $\chi(0) = 1$ .

**Definition 11** (empirical MDP for RS-KeRNS). *Let*

$$\bar{W}_h^{k+1}(x, a) = \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)).$$

In episode  $k + 1$ , RS-KeRNS uses the following estimate of the reward function

$$\check{r}_h^{k+1}(x, a) = \frac{1}{\beta + \bar{W}_h^{k+1}(x, a)} \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \tilde{r}_h^s$$

and the follow estimate of the transitions

$$\check{P}_h^{k+1}(y|x, a) = \frac{1}{\beta + \bar{W}_h^{k+1}(x, a)} \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \delta_{\bar{\zeta}_h^{s+1}(x_{h+1}^s)}(y).$$

Also, its exploration bonuses are computed as

$$\check{B}_h^{k+1}(x, a) \stackrel{\text{def}}{=} \tilde{\mathcal{O}} \left( \frac{H}{\sqrt{\beta + \bar{W}_h^{k+1}(x, a)}} + \frac{\beta H}{\beta + \bar{W}_h^{k+1}(x, a)} + L\sigma \right)$$

where the factors hidden by  $\tilde{\mathcal{O}}(\cdot)$  are the same as in Definition 5.

At step  $h$ , RS-KeRNS needs to store the quantities in Def. 11 only for the representatives  $(x, a)$  in  $\bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$  and  $y \in \bar{\mathcal{Y}}_h^{k+1}$ . We will show that, using the auxiliary quantities defined below, the values of  $\bar{W}_h^k$ ,  $\check{r}_h^k$  and  $\check{P}_h^k$  can be updated online in  $\mathcal{O}(\sum_h \bar{X}_h^k \bar{A}_h^k \bar{Y}_h^k)$  time per episode  $k$ .

**Definition 12** (auxiliary quantities for online updates). For any  $(h, x, a)$ , we define

$$\begin{aligned}\check{N}_h^{k+1}(x, a, y) &\stackrel{\text{def}}{=} \sum_{s=1}^k \chi(k-s) \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (x, a) \} \delta_{\bar{\zeta}_h^{s+1}(x_{h+1}^s)}(y) \\ \check{N}_h^{k+1}(x, a) &\stackrel{\text{def}}{=} \sum_{s=1}^k \chi(k-s) \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (x, a) \} \\ \check{S}_h^{k+1}(x, a) &\stackrel{\text{def}}{=} \sum_{s=1}^k \chi(k-s) \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (x, a) \} \check{r}_h^s.\end{aligned}$$

Notice that, if  $(x, a) \notin \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$ , the quantities above are equal to zero.

The following Lemma will be necessary in order to derive online updates.

**Lemma 17.** The empirical MDP used by *RS-KERNS* can be computed as

$$\check{r}_h^{k+1}(x, a) = \frac{\sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{S}_h^{k+1}(\bar{x}, \bar{a})}{\beta + \sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{N}_h^{k+1}(\bar{x}, \bar{a})} \quad (5)$$

$$\check{P}_h^{k+1}(y|x, a) = \frac{\sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{N}_h^{k+1}(\bar{x}, \bar{a}, y)}{\beta + \sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{N}_h^{k+1}(\bar{x}, \bar{a})} \quad (6)$$

$$\check{W}_h^{k+1}(x, a) = \sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{N}_h^{k+1}(\bar{x}, \bar{a}) \quad (7)$$

where the sums are over  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$ .

*Proof.* It is an immediate consequence of the definitions. For instance,

$$\begin{aligned}\check{W}_h^{k+1}(x, a) &= \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \\ &= \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) \sum_{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}} \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (\bar{x}, \bar{a}) \} \\ &= \sum_{(\bar{x}, \bar{a})} \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (\bar{x}, \bar{a}) \} \\ &= \sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \sum_{s=1}^k \chi(k-s) \mathbb{I} \{ \zeta_h^{s+1}(x_h^s, a_h^s) = (\bar{x}, \bar{a}) \} \\ &= \sum_{(\bar{x}, \bar{a})} \phi(\zeta_h^{k+1}(x, a), (\bar{x}, \bar{a})) \check{N}_h^{k+1}(\bar{x}, \bar{a}).\end{aligned}$$

□

---

**Algorithm 4** RS-KeRNS
 

---

- 1: **Input:** global parameters  $K, H, L, L_r, L_p, \beta, \delta, d, \sigma, \eta, W, \varepsilon_{\mathcal{X}}, \varepsilon$ .
  - 2: Initialize representative states, actions and next states:  $\bar{\mathcal{X}}_h = \emptyset, \bar{\mathcal{A}}_h = \emptyset, \bar{\mathcal{Y}}_h = \emptyset$ , for  $h \in [H]$ .
  - 3: **for** episode  $k = 1, \dots, K$  **do**
  - 4:   get initial state  $x_1^k$
  - 5:   compute  $(\tilde{Q}_h^k)_h$  using kernel backward induction on the representative sets (Alg. 5).
  - 6:   **for**  $h = 1, \dots, H$  **do**
  - 7:     execute  $a_h^k = \operatorname{argmax}_a \tilde{Q}_h^k(x_h^k, a)$ , observe reward  $\tilde{r}_h^k$  and next state  $x_{h+1}^k$
  - 8:     update representatives  $\bar{\mathcal{X}}_h, \bar{\mathcal{A}}_h, \bar{\mathcal{Y}}_h$  using  $\{x_h^k, a_h^k, x_{h+1}^k\}$  with Alg. 6
  - 9:     update model using  $x_h^k, a_h^k, x_{h+1}^k, \tilde{r}_h^k$  with Alg. 7
  - 10:   **end for**
  - 11: **end for**
- 

---

**Algorithm 5** Kernel Backward Induction on Representative States
 

---

- 1: **Input:**  $\tilde{r}_h^k(\bar{x}, \bar{a}), \tilde{P}_h^k(\bar{y}|\bar{x}, \bar{a}), \tilde{B}_h^k(\bar{x}, \bar{a})$  for all  $(\bar{x}, \bar{a}, \bar{y}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k \times \bar{\mathcal{Y}}_h^k$  and all  $h \in [H]$ .
  - 2: **Initialization:**  $\tilde{V}_{H+1}(x) = 0$  for all  $x \in \mathcal{X}$
  - 3: **for**  $h = H, \dots, 1$  **do**
  - 4:   **for**  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$  **do**
  - 5:      $\tilde{Q}_{h,\zeta}^k(\bar{x}, \bar{a}) = \tilde{r}_h^k(\bar{x}, \bar{a}) + \tilde{P}_h^k \tilde{V}_{h+1}(\bar{x}, \bar{a}) + \tilde{B}_h^k(\bar{x}, \bar{a})$
  - 6:   **end for**
  - 7:   // Interpolated  $Q$ -function. Defined, but not computed for all  $(x, a)$
  - 8:    $\tilde{Q}_h^k(x, a) = \min_{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k} \left( \tilde{Q}_{h,\zeta}^k(\bar{x}, \bar{a}) + L\rho[(x, a), (\bar{x}, \bar{a})] \right)$
  - 9:   **if**  $h > 1$  **then**
  - 10:     // Compute  $V$ -function at the next states for the stage  $h - 1$
  - 11:     **for**  $\bar{y} \in \bar{\mathcal{Y}}_{h-1}^k$  **do**
  - 12:        $\tilde{V}_h^k(\bar{y}) = \min \left( H - h + 1, \max_a \tilde{Q}_h^k(\bar{y}, a) \right)$
  - 13:     **end for**
  - 14:   **end if**
  - 15: **end for**
  - 16: **Return:**  $(\tilde{Q}_h^k)_{h \in [H]}$
- 

---

**Algorithm 6** Update Representative Sets
 

---

- 1: **Input:**  $\bar{\mathcal{X}}_h^k, \bar{\mathcal{A}}_h^k, \bar{\mathcal{Y}}_h^k, \{x_h^k, a_h^k, x_{h+1}^k\}, \varepsilon, \varepsilon_{\mathcal{X}}$ .
  - 2: **if**  $\min_{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k} \rho[(\bar{x}, \bar{a}), (x_h^k, a_h^k)] > \varepsilon$  **then**
  - 3:    $\bar{\mathcal{X}}_h^{k+1} = \bar{\mathcal{X}}_h^k \cup \{x_h^k\}, \bar{\mathcal{A}}_h^{k+1} = \bar{\mathcal{A}}_h^k \cup \{a_h^k\}$
  - 4: **else**
  - 5:    $\bar{\mathcal{X}}_h^{k+1} = \bar{\mathcal{X}}_h^k, \bar{\mathcal{A}}_h^{k+1} = \bar{\mathcal{A}}_h^k$
  - 6: **end if**
  - 7: **if**  $\min_{\bar{y} \in \bar{\mathcal{Y}}_h^k} \rho_{\mathcal{X}}(\bar{x}, x_{h+1}^k) > \varepsilon_{\mathcal{X}}$  **then**
  - 8:    $\bar{\mathcal{Y}}_h^{k+1} = \bar{\mathcal{Y}}_h^k \cup \{x_{h+1}^k\}$
  - 9: **else**
  - 10:    $\bar{\mathcal{Y}}_h^{k+1} = \bar{\mathcal{Y}}_h^k$
  - 11: **end if**
-

**Algorithm 7** Online update of RS-KERNS Model

---

```

1: Input:  $k, h, x_h^k, a_h^k, x_{h+1}^k, \tilde{r}_h^k$ .
2: // Map to representatives
3: Map  $(\tilde{x}, \tilde{a}) = \zeta_h^{k+1}(x_h^k, a_h^k)$  and  $\tilde{y} = \zeta_h^{k+1}(x_{h+1}^k)$ 
4: // Update auxiliary quantities
5:  $\tilde{N}_h^{k+1}(\tilde{x}, \tilde{a}, \tilde{y}) = 1 + \alpha_{k-s} \tilde{N}_h^k(\tilde{x}, \tilde{a}, \tilde{y})$ 
6:  $\tilde{N}_h^{k+1}(\tilde{x}, \tilde{a}) = 1 + \alpha_{k-s} \tilde{N}_h^k(\tilde{x}, \tilde{a})$ 
7:  $\tilde{S}_h^{k+1}(\tilde{x}, \tilde{a}) = \tilde{r}_h^k + \alpha_{k-s} \tilde{S}_h^k(\tilde{x}, \tilde{a})$ 
8: // Update empirical MDP
9: for  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$  do
10:   if  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$  then
11:     //  $(\bar{x}, \bar{a})$  was added before episode  $k$ 
12:      $\tilde{W}_h^{k+1}(\bar{x}, \bar{a}) = \phi((\bar{x}, \bar{a}), (\tilde{x}, \tilde{a})) + \alpha_{k-s} \tilde{W}_h^k(\bar{x}, \bar{a})$ 
13:      $\check{r}_h^{k+1}(\bar{x}, \bar{a}) = \frac{\phi((\bar{x}, \bar{a}), (\tilde{x}, \tilde{a}))}{\beta + \tilde{W}_h^{k+1}(\bar{x}, \bar{a})} \tilde{r}_h^k + \alpha_{k-s} \left( \frac{\beta + \tilde{W}_h^k(\bar{x}, \bar{a})}{\beta + \tilde{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{r}_h^k(\bar{x}, \bar{a})$ 
14:     for  $y \in \bar{\mathcal{Y}}_h^{k+1}$  do
15:        $\check{P}_h^{k+1}(y|\bar{x}, \bar{a}) = \frac{\phi((\bar{x}, \bar{a}), (\tilde{x}, \tilde{a})) \delta_{\tilde{y}}(y)}{\beta + \tilde{W}_h^{k+1}(\bar{x}, \bar{a})} + \alpha_{k-s} \left( \frac{\beta + \tilde{W}_h^k(\bar{x}, \bar{a})}{\beta + \tilde{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{P}_h^k(y|\bar{x}, \bar{a})$ 
16:     end for
17:   else
18:     //  $(\bar{x}, \bar{a})$  was added in episode  $k$ 
19:     Initialize  $\check{r}_h^{k+1}(\bar{x}, \bar{a}), \check{P}_h^{k+1}(\cdot|\bar{x}, \bar{a}), \tilde{W}_h^{k+1}(\bar{x}, \bar{a})$  using equations (5), (6) and (7)
20:   end if
21: end for

```

---

## G.2 Online updates & runtime

Assume that we observed a transition  $\{x_h^k, a_h^k, x_{h+1}^k, \tilde{r}_h^k\}$  at time  $(k, h)$ , updated the representative sets, and mapped the transition to the representatives  $(\tilde{x}, \tilde{a}, \tilde{y}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1} \times \bar{\mathcal{Y}}_h^{k+1}$ . We wish to update the estimated MDP given in Def. 11, which, at step  $h$ , are only stored for  $(x, a)$  in  $\bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$  and  $y \in \bar{\mathcal{Y}}_h^{k+1}$ .

The auxiliary quantities (Def. 12) are updated as:

$$\begin{aligned}
\tilde{N}_h^{k+1}(\tilde{x}, \tilde{a}, \tilde{y}) &= 1 + \alpha_{k-s} \tilde{N}_h^k(\tilde{x}, \tilde{a}, \tilde{y}) \\
\tilde{N}_h^{k+1}(\tilde{x}, \tilde{a}) &= 1 + \alpha_{k-s} \tilde{N}_h^k(\tilde{x}, \tilde{a}) \\
\tilde{S}_h^{k+1}(\tilde{x}, \tilde{a}) &= \tilde{r}_h^k + \alpha_{k-s} \tilde{S}_h^k(\tilde{x}, \tilde{a}).
\end{aligned}$$

We need to update  $\tilde{W}_h^k$ ,  $\check{r}_h^k$  and  $\check{P}_h^k$  for all  $(\bar{x}, \bar{a}, \bar{y}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1} \times \bar{\mathcal{Y}}_h^{k+1}$ . The update rule will depend on whether the  $(\bar{x}, \bar{a})$  is a *new* representative state-action pair (included in episode  $k$ ) or it was *visited before episode*  $k$ . These two cases are studied below.

**Case 1:**  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1}$  and  $(\bar{x}, \bar{a}) \notin \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$  This means that the representative state-action pair  $(\bar{x}, \bar{a})$  was added at time  $(k, h)$ . In this case, for all  $y \in \bar{\mathcal{Y}}_h^{k+1}$ , the quantities  $\check{r}_h^{k+1}(\bar{x}, \bar{a})$ ,  $\check{P}_h^{k+1}(y|\bar{x}, \bar{a})$  and  $\tilde{W}_h^{k+1}(\bar{x}, \bar{a})$  can be initialized using equations (5), (6) and (7). This is done in  $\mathcal{O}(\bar{\mathcal{X}}_h^{k+1} \bar{\mathcal{A}}_h^{k+1} \bar{\mathcal{Y}}_h^{k+1})$  time and can happen, at most, for one pair  $(\bar{x}, \bar{a})$ : the one that was newly added. Therefore, we have a total per-episode runtime of  $\mathcal{O}\left(\sum_{h=1}^H \bar{\mathcal{X}}_h^{k+1} \bar{\mathcal{A}}_h^{k+1} \bar{\mathcal{Y}}_h^{k+1}\right)$  taking this case into account.

**Case 2:**  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$ . This means that the representative state-action pair  $(\bar{x}, \bar{a})$  was added *before* episode  $k$ , which implies that  $\zeta_h^{k+1}(\bar{x}, \bar{a}) = \zeta_h^k(\bar{x}, \bar{a}) = (\bar{x}, \bar{a})$ . Hence,

$$\begin{aligned} \widetilde{W}_h^{k+1}(\bar{x}, \bar{a}) &= \sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(\bar{x}, \bar{a}), \zeta_h^{s+1}(x_h^s, a_h^s)) \\ &= \phi(\zeta_h^{k+1}(\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) + \sum_{s=1}^{k-1} \left( \prod_{i=1}^{k-s} \alpha_i \right) \phi(\zeta_h^k(\bar{x}, \bar{a}), \zeta_h^{s+1}(x_h^s, a_h^s)) \\ &= \phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) + \alpha_{k-s} \sum_{s=1}^{k-1} \left( \prod_{i=1}^{k-s-1} \alpha_i \right) \phi(\zeta_h^k(\bar{x}, \bar{a}), \zeta_h^{s+1}(x_h^s, a_h^s)) \\ &= \phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) + \alpha_{k-s} \widetilde{W}_h^k(\bar{x}, \bar{a}), \end{aligned}$$

where we used the convention  $\chi(0) = 1$ . This implies that, for a fixed  $(\bar{x}, \bar{a})$ , the quantity  $\widetilde{W}_h^{k+1}(\bar{x}, \bar{a})$  can be updated in  $\mathcal{O}(1)$  time, assuming that the mapping  $\zeta_h^{k+1}(x_h^k, a_h^k)$  was previously computed (this mapping is only computed *once* for all the updates, and takes  $\mathcal{O}(\bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1})$  time).

Now, notice that

$$\begin{aligned} \check{r}_h^{k+1}(\bar{x}, \bar{a}) &= \frac{\sum_{s=1}^k \chi(k-s) \phi(\zeta_h^{k+1}(\bar{x}, \bar{a}), \zeta_h^{s+1}(x_h^s, a_h^s)) \check{r}_h^s}{\beta + \widetilde{W}_h^{k+1}(\bar{x}, \bar{a})} \\ &= \frac{\phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) \check{r}_h^k}{\beta + \widetilde{W}_h^{k+1}(\bar{x}, \bar{a})} + \alpha_{k-s} \left( \frac{\beta + \widetilde{W}_h^k(\bar{x}, \bar{a})}{\beta + \widetilde{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{r}_h^k(\bar{x}, \bar{a}) \end{aligned}$$

where we used again the fact that, in this case,  $\zeta_h^{k+1}(\bar{x}, \bar{a}) = \zeta_h^k(\bar{x}, \bar{a})$ . Hence, similarly to  $\widetilde{W}_h^{k+1}(\bar{x}, \bar{a})$ , the quantity  $\check{r}_h^{k+1}(\bar{x}, \bar{a})$  can be updated in  $\mathcal{O}(1)$  time. A similar reasoning shows that  $\check{P}_h^{k+1}(y, \bar{x}, \bar{a})$  can be updated, for all  $y \in \check{\mathcal{Y}}_h^{k+1}$ , in  $\mathcal{O}(\bar{\mathcal{Y}}_h^{k+1})$  time:

$$\check{P}_h^{k+1}(y|\bar{x}, \bar{a}) = \frac{\phi((\bar{x}, \bar{a}), \zeta_h^{k+1}(x_h^k, a_h^k)) \delta_{\zeta_h^{k+1}(x_h^k, a_h^k)}(y)}{\beta + \widetilde{W}_h^{k+1}(\bar{x}, \bar{a})} + \alpha_{k-s} \left( \frac{\beta + \widetilde{W}_h^k(\bar{x}, \bar{a})}{\beta + \widetilde{W}_h^{k+1}(\bar{x}, \bar{a})} \right) \check{P}_h^k(y|\bar{x}, \bar{a}).$$

**Summary** Every time a new transition is observed at time  $(k, h)$ , the estimators for all  $(x, a, y) \in \bar{\mathcal{X}}_h^{k+1} \times \bar{\mathcal{A}}_h^{k+1} \times \check{\mathcal{Y}}_h^{k+1}$  must be updated. For a given representative  $(x, a)$ , the updates can be done in  $\mathcal{O}(\bar{\mathcal{Y}}_h^{k+1})$  time if it has been observed before episode  $k$  (case 2). This results in a total runtime, per episode, of  $\mathcal{O}(\sum_h \bar{\mathcal{X}}_h^{k+1} \bar{\mathcal{A}}_h^{k+1} \bar{\mathcal{Y}}_h^{k+1})$  for all the representatives observed before episode  $k$ . If the representative  $(x, a)$  has not been observed before episode  $k$  (case 1), the updates require  $\mathcal{O}(\bar{\mathcal{X}}_h^{k+1} \bar{\mathcal{A}}_h^{k+1} \bar{\mathcal{Y}}_h^{k+1})$  time, and this can happen, at most, for one state-action pair at each time  $(k, h)$ . Hence, the total runtime required for the updates is  $\mathcal{O}(\sum_h \bar{\mathcal{X}}_h^{k+1} \bar{\mathcal{A}}_h^{k+1} \bar{\mathcal{Y}}_h^{k+1})$  per episode.

### G.3 Regret analysis

The regret analysis of **RS-KeRNS** is based on the following result, which is a corollary of Lemma 25, and is used to bound the bias introduced by using representative states.

**Corollary 10.** Let  $\chi_{(\eta, W)} : \mathbb{N} \rightarrow [0, 1]$ , consider the following kernel

$$\Gamma(t, u, v) = \chi_{(\eta, W)}(t) \exp\left(-\rho[u, v]^2 / (2\sigma^2)\right)$$

where  $u, v \in \mathcal{X} \times \mathcal{A}$ . For  $s < k$ , let

$$w_h^{k,s}(x, a) = \Gamma(k - s - 1, (x, a), (x_h^s, a_h^s)), \quad w_{h,\zeta}^{k,s}(x, a) \stackrel{\text{def}}{=} \Gamma(k - s - 1, \zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))$$

and consider the functions

$$\begin{aligned} g_1(x, a) &= \frac{\sum_{s=1}^{k-1} w_h^{k,s}(x, a) Y_s}{\beta + \sum_{s=1}^{k-1} w_h^{k,s}(x, a)}, & g_1^\zeta(x, a) &= \frac{\sum_{s=1}^{k-1} w_{h,\zeta}^{k,s}(x, a) Y_s}{\beta + \sum_{s=1}^{k-1} w_{h,\zeta}^{k,s}(x, a)}, \\ g_2(x, a) &= \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} w_h^{k,s}(x, a)}}, & g_2^\zeta(x, a) &= \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} w_{h,\zeta}^{k,s}(x, a)}}, \\ g_3(x, a) &= \frac{1}{\beta + \sum_{s=1}^{k-1} w_h^{k,s}(x, a)}, & g_3^\zeta(x, a) &= \frac{1}{\beta + \sum_{s=1}^{k-1} w_{h,\zeta}^{k,s}(x, a)}. \end{aligned}$$

where  $(Y_s)_{s=1}^{k-1}$  is an arbitrary sequence. Then,  $g_1$ ,  $g_2$  and  $g_3$  are Lipschitz continuous, whose Lipschitz constants are bounded by  $L_1$ ,  $L_2$  and  $L_3$ , respectively, with

$$L_1 = \frac{4 \max_s |Y_s|}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right), \quad L_2 = \frac{1 + \sqrt{\log^+(k/\beta)}}{2\beta^{1/2}\sigma}, \quad L_3 = \frac{1 + \sqrt{\log^+(k/\beta)}}{\beta\sigma}$$

Furthermore, for any  $(x, a)$  and for  $i \in \{1, 2, 3\}$ ,

$$\left|g_i^\zeta(x, a) - g_i(\zeta_h^k(x, a))\right| \leq L_i \max_s \rho[(x_h^s, a_h^s), \zeta_h^k(x_h^s, a_h^s)].$$

*Proof.* First, let's prove that  $g_1$ ,  $g_2$  and  $g_3$  are Lipschitz continuous. From Lemma 25, taking  $z = (\rho[(x, a), (x_h^s, a_h^s)])_{s=1}^{k-1}$  and  $y = (\rho[(x', a'), (x_h^s, a_h^s)])_{s=1}^{k-1}$  we have

$$|g_1(x, a) - g_1(x', a')| \leq \frac{4 \max_s |Y_s|}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_{s \in [k-1]} |z_s - y_s|.$$

By the triangle inequality, for all  $s$ ,

$$|z_s - y_s| = |\rho[(x, a), (x_h^s, a_h^s)] - \rho[(x', a'), (x_h^s, a_h^s)]| \leq \rho[(x, a), (x', a')]$$

which implies

$$|g_1(x, a) - g_1(x', a')| \leq \frac{4 \max_s |Y_s|}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \rho[(x, a), (x', a')].$$

giving the Lipschitz constant of  $g_1$ . The Lipschitz constants of  $g_2$  and  $g_3$  follow similarly from Lemma 25.

Now, let's bound the differences  $\left|g_i^\zeta(x, a) - g_i(\zeta_h^k(x, a))\right|$ . Let  $\chi(s) = \chi_{(\eta, W)}(k - s - 1)$ . For  $i = 1$ , and applying

again Lemma 25, we obtain

$$\begin{aligned}
 & \left| g_1^\zeta(x, a) - g_1(\zeta_h^k(x, a)) \right| \\
 &= \left| \frac{\sum_{s=1}^{k-1} \chi(s) \exp\left(-\frac{\rho[\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)]^2}{2\sigma^2}\right) Y_s}{\beta + \sum_{s=1}^{k-1} \chi(s) \exp\left(-\frac{\rho[\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)]^2}{2\sigma^2}\right)} - \frac{\sum_{s=1}^{k-1} \chi(s) \exp\left(-\frac{\rho[\zeta_h^k(x, a), (x_h^s, a_h^s)]^2}{2\sigma^2}\right) Y_s}{\beta + \sum_{s=1}^{k-1} \chi(s) \exp\left(-\frac{\rho[\zeta_h^k(x, a), (x_h^s, a_h^s)]^2}{2\sigma^2}\right)} \right| \\
 &\leq \frac{4 \max_s |Y_s|}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_{s \in [k-1]} \left| \rho[\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)] - \rho[\zeta_h^k(x, a), (x_h^s, a_h^s)] \right| \\
 &\leq \frac{4 \max_s |Y_s|}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_{s \in [k-1]} \rho[\zeta_h^k(x_h^s, a_h^s), (x_h^s, a_h^s)],
 \end{aligned}$$

where, in the last line, we used the triangle inequality. The proof for  $i \in \{2, 3\}$  also follow from Lemma 25.  $\square$

We use Corollary 10 to bound the difference between the estimators and the bonuses of **KeRNS**,  $(\hat{r}_h^k, \hat{P}_h^k, \mathbf{B}_h^k)$ , and the ones of **RS-KeRNS**,  $(\check{r}_h^k, \check{P}_h^k, \check{\mathbf{B}}_h^k)$ .

**Lemma 18.** *Let  $V$  be an arbitrary  $L$ -Lipschitz function bounded by  $H$ . Then, for any  $(x, a)$ ,*

$$\begin{aligned}
 \left| (\hat{P}_h^k - \check{P}_h^k) V(x, a) \right| &\leq \frac{4H}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) (\rho[(x, a), \zeta_h^k(x, a)] + \varepsilon) \\
 &\quad + 3L\varepsilon_{\mathcal{X}} + 8H \left(1 + \sqrt{\log^+(k/\beta)}\right) \frac{\varepsilon}{\sigma}
 \end{aligned}$$

*Proof.* To simplify the notations, let  $f(s) \stackrel{\text{def}}{=} \chi(k-s-1)$  for  $s \in [k-1]$ . We have

$$\begin{aligned}
 & \left| (\hat{P}_h^k - \check{P}_h^k) V(x, a) \right| \\
 &= \left| \frac{\sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s)) V(x_{h+1}^s)}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) V(\bar{\zeta}_h^{s+1}(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))} \right| \\
 &\leq \textcircled{1} + \textcircled{2}
 \end{aligned}$$

where  $\textcircled{1}$  and  $\textcircled{2}$  are defined and bounded below. First,

$$\begin{aligned}
 \textcircled{1} &= \left| \frac{\sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s)) V(x_{h+1}^s)}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))} \right| \\
 &\leq \left| \frac{\sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s)) V(x_{h+1}^s)}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} \right| \\
 &\quad + \left| \frac{\sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), (x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), (x_h^s, a_h^s))} \right| \\
 &\quad + \left| \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), (x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), (x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))} \right| \\
 &\leq L \max_{s \in [k-1]} \rho_{\mathcal{X}}(x_{h+1}^s, \bar{\zeta}_h^k(x_{h+1}^s)) \\
 &\quad + \frac{4H}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \rho[(x, a), \zeta_h^k(x, a)] \\
 &\quad + \frac{4H}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_s \rho[(x_h^s, a_h^s), \zeta_h^k(x_h^s, a_h^s)]
 \end{aligned}$$

by using the fact that  $V$  is  $L$ -Lipschitz and Corollary 10.



To bound ②, let  $z_s = \rho [\zeta^k(x, a), \zeta^{s+1}(x_h^s, a_h^s)]$  and  $y_s = \rho [\zeta^k(x, a), \zeta^k(x_h^s, a_h^s)]$ , for  $s \in [k-1]$ . Using again the fact that  $V$  is  $L$ -Lipschitz and Corollary 10, we obtain

$$\begin{aligned} \textcircled{2} &= \left| \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s)) V(\bar{\zeta}_h^k(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))} - \frac{\sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s)) V(\bar{\zeta}_h^{s+1}(x_{h+1}^s))}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))} \right| \\ &\leq L \max_s \rho_{\mathcal{X}}(\bar{\zeta}^{s+1}(x_{h+1}^s), \bar{\zeta}^k(x_{h+1}^s)) + \frac{4H}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_s |z_s - y_s| \\ &\leq 2L\varepsilon_{\mathcal{X}} + \frac{4H}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) \max_s |\rho[\zeta^k(x_h^s, a_h^s), \zeta^{s+1}(x_h^s, a_h^s)]| \\ &\leq 2L\varepsilon_{\mathcal{X}} + 8H \left(1 + \sqrt{\log^+(k/\beta)}\right) \frac{\varepsilon}{\sigma}. \end{aligned}$$

By the construction of the algorithm,  $\rho_{\mathcal{X}}(x_{h+1}^s, \bar{\zeta}_h^k(x_{h+1}^s)) \leq \varepsilon_{\mathcal{X}}$ ,  $\rho[(x_h^s, a_h^s), \zeta_h^k(x_h^s, a_h^s)] \leq \varepsilon$ ,  $\rho_{\mathcal{X}}(\bar{\zeta}^{s+1}(x_{h+1}^s), \bar{\zeta}^k(x_{h+1}^s)) \leq 2\varepsilon_{\mathcal{X}}$  and  $\rho[\zeta^k(x_h^s, a_h^s), \zeta^{s+1}(x_h^s, a_h^s)] \leq 2\varepsilon$ , which concludes the proof.  $\square$

**Lemma 19.** For any  $(x, a)$ ,

$$\begin{aligned} |\hat{r}_h^k(x, a) - \check{r}_h^k(x, a)| &\leq \frac{4}{\sigma} \left(1 + \sqrt{\log^+(k/\beta)}\right) (\rho[(x, a), \zeta_h^k(x, a)] + \varepsilon) \\ &\quad + 8 \left(1 + \sqrt{\log^+(k/\beta)}\right) \frac{\varepsilon}{\sigma} \end{aligned}$$

*Proof.* It follows from a similar decomposition as in the proof of Lemma 18, and from Corollary 10.  $\square$

**Lemma 20.** For any  $(x, a)$ , we have,

$$\left| \mathbb{B}_h^k(x, a) - \check{\mathbb{B}}_h^k(x, a) \right| \lesssim \frac{H}{\sigma} \left(1 + \frac{1}{2\beta^{1/2}}\right) (\rho[(x, a), \zeta_h^k(x, a)] + \varepsilon) + H \left(1 + \frac{1}{\beta^{1/2}}\right) \frac{\varepsilon}{\sigma}$$

*Proof.* To simplify the notations, let  $f(s) \stackrel{\text{def}}{=} \chi(k-s-1)$  for  $s \in [k-1]$ . Using definitions 5, 10, and 11, we have

$$\begin{aligned} &\left| \mathbb{B}_h^k(x, a) - \check{\mathbb{B}}_h^k(x, a) \right| \\ &\lesssim H \left| \underbrace{\sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))}} - \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))}}}_{\textcircled{1}} \right| \\ &\quad + \beta H \left| \underbrace{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s) \phi((x, a), (x_h^s, a_h^s))} - \frac{1}{\beta + \sum_{s=1}^{k-1} f(s) \phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))}}}_{\textcircled{2}} \right|. \end{aligned}$$

Using Corollary 10, we obtain

$$\begin{aligned}
 \textcircled{1} &\leq \left| \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi((x, a), (x_h^s, a_h^s))}} - \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi(\zeta_h^k(x, a), (x_h^s, a_h^s))}} \right| \\
 &+ \left| \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi(\zeta_h^k(x, a), (x_h^s, a_h^s))}} - \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))}} \right| \\
 &+ \left| \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi(\zeta_h^k(x, a), \zeta_h^k(x_h^s, a_h^s))}} - \sqrt{\frac{1}{\beta + \sum_{s=1}^{k-1} f(s)\phi(\zeta_h^k(x, a), \zeta_h^{s+1}(x_h^s, a_h^s))}} \right| \\
 &\leq \left( \frac{1 + \sqrt{\log^+(k/\beta)}}{2\beta^{1/2}\sigma} \right) \left( \rho[(x, a), \zeta_h^k(x, a)] + \max_{s \in [k-1]} \rho[(x_h^s, a_h^s), \zeta_h^k(x, a)] \right. \\
 &\quad \left. + \max_{s \in [k-1]} \rho[\zeta_h^k(x_h^s, a_h^s), \zeta_h^{s+1}(x, a)] \right).
 \end{aligned}$$

Similarly, Corollary 10 yields

$$\begin{aligned}
 \textcircled{2} &\leq \left( \frac{1 + \sqrt{\log^+(k/\beta)}}{\beta\sigma} \right) \left( \rho[(x, a), \zeta_h^k(x, a)] + \max_{s \in [k-1]} \rho[(x_h^s, a_h^s), \zeta_h^k(x, a)] \right. \\
 &\quad \left. + \max_{s \in [k-1]} \rho[\zeta_h^k(x_h^s, a_h^s), \zeta_h^{s+1}(x, a)] \right).
 \end{aligned}$$

By the construction of the algorithm,  $\rho[(x_h^s, a_h^s), \zeta_h^k(x, a)] \leq \varepsilon$  and  $\rho[\zeta_h^k(x_h^s, a_h^s), \zeta_h^{s+1}(x_h^s, a_h^s)] \leq 2\varepsilon$ , which concludes the proof.  $\square$

**Lemma 21.** Let  $\tilde{Q}_h^k$  and  $\tilde{Q}_{h,\zeta}^k$  be the  $Q$ -functions defined in Algorithm 5. Then,

$$\begin{aligned}
 \tilde{Q}_h^k(x, a) &\stackrel{\text{def}}{=} \min_{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k} \left( \tilde{Q}_{h,\zeta}^k(\bar{x}, \bar{a}) + L\rho[(x, a), (\bar{x}, \bar{a})] \right) \\
 &= \min_{s \in [k-1]} \left( \tilde{Q}_{h,\zeta}^k(x_h^s, a_h^s) + L\rho[(x, a), \zeta_h^k(x_h^s, a_h^s)] \right).
 \end{aligned}$$

*Proof.* Notice that, although  $\tilde{Q}_{h,\zeta}^k$  is only computed for the representative state-action pairs, it is defined for any  $(x, a)$  as

$$\tilde{Q}_{h,\zeta}^k(x, a) = \check{r}_h^k(x, a) + \check{P}_h^k \check{V}_{h+1}(x, a) + \check{B}_h^k(x, a).$$

We claim that

$$\underbrace{\{\zeta_h^k(x_h^s, a_h^s) : s \in [k-1]\}}_A = \underbrace{\bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k}_B.$$

First,  $A \subset B$ , since  $\forall (s, h)$ , we have  $\zeta_h^k(x_h^s, a_h^s) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$ . Second, for any  $(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k$ , there exists  $(s, h)$  such that  $(\bar{x}, \bar{a}) = (x_h^s, a_h^s) = \zeta_h^{s+1}(x_h^s, a_h^s) = \zeta_h^k(x_h^s, a_h^s) \in A$ , which implies that  $B \subset A$ .

Together with the fact that  $\tilde{Q}_{h,\zeta}^k(x_h^s, a_h^s) = \tilde{Q}_{h,\zeta}^k(\zeta_h^k(x_h^s, a_h^s))$ , which holds by Definition 11 and Algorithm 5, this concludes the proof.  $\square$

**Lemma 22.** The difference between the  $Q$ -values computed by *KeRNS* and *RS-KeRNS* are bounded as follows

$$\sup_{x,a} \left| Q_h^k(x, a) - \tilde{Q}_h^k(x, a) \right| \lesssim \left( L(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma} H \right) (H - h + 1).$$

*Proof.* We proceed by induction on  $h$ . Let

$$\epsilon_h \stackrel{\text{def}}{=} \left( L(\epsilon + \epsilon_{\mathcal{X}}) + \frac{\epsilon}{\sigma} H \right) (H - h + 1).$$

For  $h = H + 1$ ,  $\epsilon_{H+1} = 0$ ,  $Q_{H+1}^k = \check{Q}_{H+1}^k = 0$  and the claim holds.

Now, assume that the claim is true for  $h + 1$ . In this case, we have, for any  $x$ ,

$$\begin{aligned} & \left| V_{h+1}^k(x) - \check{V}_{h+1}^k(x) \right| \\ &= \left| \min \left( H - h, \max_a Q_{h+1}^k(x, a) \right) - \min \left( H - h, \max_a \check{Q}_{h+1}^k(x, a) \right) \right| \\ &\leq \left| \max_a Q_{h+1}^k(x, a) - \max_a \check{Q}_{h+1}^k(x, a) \right| \\ &\leq \max_a \left| Q_{h+1}^k(x, a) - \check{Q}_{h+1}^k(x, a) \right| \leq \epsilon_{h+1}, \end{aligned}$$

where we used induction hypothesis and the fact that, for any  $a, b, c \in \mathbb{R}$ , we have  $|\min(a, b) - \min(a, c)| \leq |b - c|$  (Fact 1).

For any  $(x_h^s, a_h^s)$  with  $s \in [k - 1]$ , we have

$$\begin{aligned} & \left| \tilde{Q}_h^k(x_h^s, a_h^s) - \check{Q}_{h, \zeta}^k(x_h^s, a_h^s) \right| \\ &\leq \left| \check{r}_h^k(x_h^s, a_h^s) - \check{r}_{h, \zeta}^k(x_h^s, a_h^s) \right| + \left| \mathbb{B}_h^k(x_h^s, a_h^s) - \check{\mathbb{B}}_h^k(x_h^s, a_h^s) \right| \\ &\quad + \left| \left( \hat{P}_h^k - \check{P}_h^k \right) V_{h+1}^k(x_h^s, a_h^s) + \check{P}_h^k \left( V_{h+1}^k - \check{V}_{h+1}^k \right) (x_h^s, a_h^s) \right| \\ &\lesssim \epsilon_{h+1} + \frac{H}{\sigma} \left( \rho \left[ (x_h^s, a_h^s), \zeta_h^k(x_h^s, a_h^s) \right] + \epsilon \right) + L\epsilon_{\mathcal{X}} + \frac{\epsilon}{\sigma} H \end{aligned}$$

where, in the last line, we used the induction hypothesis and lemmas 18, 19 and 20.

By the construction of **RS-KerNS**, we have  $\max_{s' \in [k-1]} \rho \left[ (x_h^{s'}, a_h^{s'}), \zeta_h^k(x_h^{s'}, a_h^{s'}) \right] \leq \epsilon$ . Consequently,

$$\left| \tilde{Q}_h^k(x_h^s, a_h^s) - \check{Q}_{h, \zeta}^k(x_h^s, a_h^s) \right| \lesssim \epsilon_{h+1} + L\epsilon_{\mathcal{X}} + \frac{\epsilon}{\sigma} H.$$

Now, take an arbitrary  $(x, a)$ . By Lemma 21,

$$\begin{aligned} \check{Q}_h^k(x, a) &= \min_{(\bar{x}, \bar{a}) \in \check{\mathcal{X}}_h^k \times \check{\mathcal{A}}_h^k} \left( \check{Q}_{h, \zeta}^k(\bar{x}, \bar{a}) + L\rho \left[ (x, a), (\bar{x}, \bar{a}) \right] \right) \\ &= \min_{s \in [k-1]} \left( \check{Q}_{h, \zeta}^k(x_h^s, a_h^s) + L\rho \left[ (x, a), \zeta_h^k(x_h^s, a_h^s) \right] \right) \end{aligned}$$

and, by definition,

$$Q_h^k(x, a) = \min_{s \in [k-1]} \left( \tilde{Q}_h^k(x_h^s, a_h^s) + L\rho \left[ (x, a), (x_h^s, a_h^s) \right] \right),$$

we obtain, for any  $(x, a)$ ,

$$\begin{aligned} & \left| Q_h^k(x, a) - \check{Q}_h^k(x, a) \right| \\ &\lesssim \min_{s \in [k-1]} \left| \tilde{Q}_h^k(x_h^s, a_h^s) - \check{Q}_{h, \zeta}^k(x_h^s, a_h^s) \right| + L \min_{s \in [k-1]} \left| \rho \left[ (x, a), (x_h^s, a_h^s) \right] - \rho \left[ (x, a), \zeta_h^k(x_h^s, a_h^s) \right] \right| \\ &\leq \epsilon_{h+1} + L\epsilon_{\mathcal{X}} + \frac{\epsilon}{\sigma} H + L \max_{s \in [k-1]} \rho \left[ (x_h^s, a_h^s), \zeta_h^k(x_h^s, a_h^s) \right] \\ &\leq \epsilon_{h+1} + L(\epsilon + \epsilon_{\mathcal{X}}) + \frac{\epsilon}{\sigma} H = \epsilon_h. \end{aligned}$$

which concludes the proof.  $\square$

**Theorem 5** (UCRL-type regret bound for **RS-KeRNS**). *With probability at least  $1 - \delta$ , the regret of **RS-KeRNS** is bounded as follows*

$$\mathcal{R}^{\text{RS-KeRNS}}(K) \lesssim \mathcal{R}_1^{\text{KeRNS}}(K) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3.$$

where  $\mathcal{R}_1^{\text{KeRNS}}(K)$  is the UCRL-type regret bound given in Theorem 3 for **KeRNS**.

*Proof.* Let  $\check{\pi}_k$  be the policy followed by **RS-KeRNS** in episode  $k$  and let  $\check{\delta}_h^k \stackrel{\text{def}}{=} \check{V}_h^k(x_h^k) - V_{k,h}^{\check{\pi}_k}(x_h^k)$

**Regret decomposition** Consider the following decomposition, also used in the proof of Lemma 11:

$$\begin{aligned} \check{\delta}_h^k &= \check{V}_h^k(x_h^k) - V_{k,h}^{\check{\pi}_k}(x_h^k) \\ &\leq \check{Q}_h^k(x_h^k, a_h^k) - Q_{k,h}^{\check{\pi}_k}(x_h^k, a_h^k) \\ &\leq \check{Q}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - Q_{k,h}^{\check{\pi}_k}(x_h^k, a_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)], \quad \text{since } \check{Q}_h^k \text{ is } L\text{-Lipschitz} \\ &\leq \check{Q}_{h,\zeta}^k(\tilde{x}_h^k, \tilde{a}_h^k) - Q_{k,h}^{\check{\pi}_k}(x_h^k, a_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)], \quad \text{since } \check{Q}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \leq \check{Q}_{h,\zeta}^k(\tilde{x}_h^k, \tilde{a}_h^k) \\ &= \check{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k) + \check{P}_h^k V_{h+1}^k(\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k V_{k,h+1}^{\check{\pi}_k}(x_h^k, a_h^k) + \check{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)] \\ &= \underbrace{\check{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k)}_{\text{(A)}} + \underbrace{[\check{P}_h^k - P_h^k] V_{k,h+1}^*(\tilde{x}_h^k, \tilde{a}_h^k)}_{\text{(B)}} + \underbrace{[\check{P}_h^k - P_h^k] (\check{V}_{h+1}^k - V_{k,h+1}^*) (\tilde{x}_h^k, \tilde{a}_h^k)}_{\text{(C)}} \\ &\quad + \underbrace{P_h^k \check{V}_{h+1}^k(\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k V_{k,h+1}^{\check{\pi}_k}(x_h^k, a_h^k)}_{\text{(D)}} + \check{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + 2L\rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)]. \end{aligned}$$

We will use the following results, which are a consequence of Lemmas 18, 19 and 20 and the fact that  $\rho [(\tilde{x}_h^k, \tilde{a}_h^k), \zeta(\tilde{x}_h^k, \tilde{a}_h^k)] \leq \varepsilon$ :

$$|\check{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k)| \lesssim \frac{4}{\sigma} (\rho [(\tilde{x}_h^k, \tilde{a}_h^k), \zeta(\tilde{x}_h^k, \tilde{a}_h^k)] + \varepsilon) + \frac{8\varepsilon}{\sigma} \lesssim \frac{\varepsilon}{\sigma}.$$

$$\begin{aligned} \check{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) &= B_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \check{B}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - B_h^k(\tilde{x}_h^k, \tilde{a}_h^k) \\ &\lesssim B_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \frac{H}{\sigma} (\rho [(\tilde{x}_h^k, \tilde{a}_h^k), \zeta(\tilde{x}_h^k, \tilde{a}_h^k)] + \varepsilon) + \frac{H\varepsilon}{\sigma} \\ &\lesssim B_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \frac{\varepsilon}{\sigma}H. \end{aligned}$$

and, for any function  $f$  that is  $L$ -Lipschitz and bounded by  $H$ ,

$$\left| (\check{P}_h^k - P_h^k) f(\tilde{x}_h^k, \tilde{a}_h^k) \right| \lesssim L\varepsilon_{\mathcal{X}} + \frac{4H}{\sigma} (\rho [(\tilde{x}_h^k, \tilde{a}_h^k), \zeta(\tilde{x}_h^k, \tilde{a}_h^k)] + \varepsilon) + \frac{H\varepsilon}{\sigma} \lesssim L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma}H.$$

Now, we bound each term (A) – (D).

Term (A): by Lemma 19, the definition of  $\mathcal{G}$  and Corollary 2, we have

$$\begin{aligned} \text{(A)} &= \check{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k) \\ &\leq \frac{\varepsilon}{\sigma} + \hat{r}_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + r_h^k(\tilde{x}_h^k, \tilde{a}_h^k) - r_h^k(x_h^k, a_h^k) \\ &\leq \frac{\varepsilon}{\sigma} + r B_h^k(\tilde{x}_h^k, \tilde{a}_h^k) + \text{bias}_r(k, h) + L_r \rho [(\tilde{x}_h^k, \tilde{a}_h^k), (x_h^k, a_h^k)]. \end{aligned}$$

Term **(B)**:

$$\begin{aligned}
 \mathbf{(B)} &= \left[ \check{P}_h^k - \hat{P}_h^k \right] \mathbf{V}_{k,h+1}^* (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \hat{P}_h^k - P_h^k \right] \mathbf{V}_{k,h+1}^* (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\lesssim L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H + \left[ \hat{P}_h^k - \bar{P}_h^k \right] \mathbf{V}_{k,h+1}^* (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \bar{P}_h^k - P_h^k \right] \mathbf{V}_{k,h+1}^* (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\leq L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H + {}^p \mathbf{B}_h^k (\tilde{x}_h^k, \tilde{a}_h^k) + \mathbf{bias}_p(k, h)
 \end{aligned}$$

Term **(C)**: Using Corollary 2, we obtain

$$\begin{aligned}
 \mathbf{(C)} &= \left[ \check{P}_h^k - P_h^k \right] \left( \check{V}_{h+1}^k - \mathbf{V}_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &= \left[ \check{P}_h^k - \hat{P}_h^k \right] \left( \check{V}_{h+1}^k - \mathbf{V}_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \hat{P}_h^k - P_h^k \right] \left( \check{V}_{h+1}^k - \mathbf{V}_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\lesssim L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H + \left[ \hat{P}_h^k - \bar{P}_h^k \right] \left( \check{V}_{h+1}^k - \mathbf{V}_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \mathbf{bias}_p(k, h) \\
 &\lesssim L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H + \sqrt{\frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \mathbf{bias}_p(k, h)
 \end{aligned}$$

by the definition of  $\mathcal{G}$ .

Term **(D)**: From Assumption 2, for any  $L$ -Lipschitz function, the mapping  $(x, a) \mapsto P_h^k f(x, a)$  is  $L_p L$ -Lipschitz. Consequently,

$$\begin{aligned}
 \mathbf{(D)} &= P_h^k \check{V}_{h+1}^k (\tilde{x}_h^k, \tilde{a}_h^k) - P_h^k \check{V}_{k,h+1}^{\pi_k} (x_h^k, a_h^k) \\
 &\leq P_h^k \check{V}_{h+1}^k (x_h^k, a_h^k) - P_h^k \check{V}_{k,h+1}^{\pi_k} (x_h^k, a_h^k) + L_p L \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &= P_h^k \left( \check{V}_{h+1}^k - \check{V}_{k,h+1}^{\pi_k} \right) (x_h^k, a_h^k) + L_p L \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \\
 &= \check{\delta}_{h+1}^k + \xi_{h+1}^k + L_p L \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)].
 \end{aligned}$$

where

$$\xi_{h+1}^k \stackrel{\text{def}}{=} P_h^k \left( \check{V}_{h+1}^k - \check{V}_{k,h+1}^{\pi_k} \right) (x_h^k, a_h^k) - \check{\delta}_{h+1}^k$$

is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .

Putting together the bounds for **(A)**-**(D)** and using the definition of the bonuses  $\mathbf{B}_h^k$ , we obtain

$$\begin{aligned}
 \check{\delta}_h^k &\lesssim \check{\delta}_{h+1}^k + \xi_{h+1}^k + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] + \sqrt{\frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)}} \\
 &\quad + \frac{\beta H}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma + \mathbf{bias}(k, h') + L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H
 \end{aligned}$$

where the constant in front of  $\check{\delta}_{h+1}^k$  is *exact* (i.e., not omitted by  $\lesssim$ ).

Now, we follow the same arguments as in the proof of Lemma 11. Consider the event  $E_h^k \stackrel{\text{def}}{=} \{\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$ . The inequality above gives us

$$\begin{aligned}
 \mathbb{I} \{E_h^k\} \check{\delta}_h^k &\lesssim \mathbb{I} \{E_h^k\} \check{\delta}_{h+1}^k + \mathbb{I} \{E_h^k\} \xi_{h+1}^k + \mathbb{I} \{E_h^k\} \sqrt{\frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)}} \\
 &\quad + \mathbb{I} \{E_h^k\} \frac{\beta H}{\mathbf{C}_h^k (\tilde{x}_h^k, \tilde{a}_h^k)} + 3L\sigma + \mathbf{bias}(k, h') + L\varepsilon_{\mathcal{X}} + \frac{\varepsilon}{\sigma} H.
 \end{aligned} \tag{8}$$

Using Lemmas 10 and 22, we upper bound  $\mathbb{I}\{E_h^k\} \check{\delta}_{h+1}^k$  in terms of  $\check{\delta}_{h+1}^k$ :

$$\begin{aligned}
 \mathbb{I}\{E_h^k\} \check{\delta}_{h+1}^k &= \mathbb{I}\{E_h^k\} \left( \check{V}_{h+1}^k(x_{h+1}^k) - V_{k,h}^{\check{\pi}_k}(x_{h+1}^k) \right) \\
 &\lesssim \mathbb{I}\{E_h^k\} \left( HL(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma} H^2 + V_{h+1}^k(x_{h+1}^k) - V_{k,h+1}^{\check{\pi}_k}(x_{h+1}^k) \right) \\
 &\leq HL(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma} H^2 + \mathbb{I}\{E_h^k\} \underbrace{\left( V_{h+1}^k(x_{h+1}^k) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) - V_{k,h+1}^{\check{\pi}_k}(x_{h+1}^k) \right)}_{\geq 0 \text{ by Lemma 10}} \\
 &\leq HL(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma} H^2 + V_{h+1}^k(x_{h+1}^k) - V_{k,h+1}^{\check{\pi}_k}(x_{h+1}^k) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \\
 &\leq 2HL(\varepsilon + \varepsilon_{\mathcal{X}}) + 2\frac{\varepsilon}{\sigma} H^2 + \check{V}_{h+1}^k(x_{h+1}^k) - V_{k,h+1}^{\check{\pi}_k}(x_{h+1}^k) + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \\
 &= \check{\delta}_{h+1}^k + 2HL(\varepsilon + \varepsilon_{\mathcal{X}}) + 2\frac{\varepsilon}{\sigma} H^2 + \sum_{h'=h+1}^H \mathbf{bias}(k, h) \tag{9}
 \end{aligned}$$

Let  $\overline{E}_h^k$  be the complement of  $E_h^k$ . Using the inequality above combined with (8), and the fact that  $\check{\delta}_h^k \leq H$ , we obtain

$$\begin{aligned}
 \check{\delta}_h^k &= \mathbb{I}\{\overline{E}_h^k\} \check{\delta}_h^k + \mathbb{I}\{E_h^k\} \check{\delta}_h^k \\
 &\leq H \mathbb{I}\{\overline{E}_h^k\} + \mathbb{I}\{E_h^k\} \check{\delta}_h^k \\
 &\lesssim H \mathbb{I}\{\overline{E}_h^k\} + \check{\delta}_{h+1}^k + HL(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma} H^2 + \sum_{h'=h}^H \mathbf{bias}(k, h) + \mathbb{I}\{E_h^k\} \xi_{h+1}^k \\
 &\quad + \mathbb{I}\{E_h^k\} \sqrt{\frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \mathbb{I}\{E_h^k\} \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} + L\sigma
 \end{aligned} \tag{10}$$

Consequently,

$$\begin{aligned}
 \sum_{k=1}^K \check{\delta}_1^k &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H \sqrt{|\mathcal{C}'_{\sigma}|}}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I}\{E_h^k\} + H \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{\overline{E}_h^k\} \\
 &\quad + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \xi_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + L(\varepsilon + \varepsilon_{\mathcal{X}}) K H^2 + \frac{\varepsilon}{\sigma} K H^3
 \end{aligned}$$

Now, as in the proof of Lemma 11, we show that

$$H \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{\overline{E}_h^k\} \leq H^2 |\mathcal{C}_{\sigma}|$$

and, from lemmas 10 and 22, we have:

$$\begin{aligned}
 \mathcal{R}^{\text{RS-KeRNS}}(K) &= \sum_{k=1}^K \left( V_{k,1}^*(x_1^k) - V_{k,h}^{\tilde{\pi}_k}(x_1^k) \right) \\
 &\leq \sum_{k=1}^K \left( V_1^k(x_1^k) - V_{k,h}^{\tilde{\pi}_k}(x_1^k) \right) + \sum_{k=1}^K \sum_{h=1}^H \text{bias}(k, h) \\
 &\lesssim \sum_{k=1}^K \left( \check{V}_1^k(x_1^k) - V_{k,h}^{\tilde{\pi}_k}(x_1^k) \right) + \sum_{k=1}^K \sum_{h=1}^H \text{bias}(k, h) + LKH(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma}KH^2 \\
 &= \sum_{k=1}^K \check{\delta}_1^k + \sum_{k=1}^K \sum_{h=1}^H \text{bias}(k, h) + LKH(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma}KH^2 \\
 &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H\sqrt{|\mathcal{C}'_{\sigma}|}}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{\beta H}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I}\{E_h^k\} + H^2|\mathcal{C}_{\sigma}| \\
 &\quad + \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}\{E_h^k\} \xi_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \text{bias}(k, h) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3.
 \end{aligned}$$

Recall the definition  $E_h^k \stackrel{\text{def}}{=} \{\rho[(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma\}$  and the fact that  $\tilde{\xi}_{h+1}^k \stackrel{\text{def}}{=} \mathbb{I}\{E_h^k\} \xi_{h+1}^k$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ . Then, as in the proof of Theorem 3, we obtain

$$\mathcal{R}^{\text{RS-KeRNS}}(K) \lesssim \mathcal{R}_1^{\text{KeRNS}}(K) + LKH^2(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma}KH^3$$

with probability at least  $1 - \delta$ . □

**Theorem 6** (UCBVI-type regret bound for RS-KeRNS). *With probability at least  $1 - \delta$ , the regret of RS-KeRNS is bounded as follows*

$$\mathcal{R}^{\text{RS-KeRNS}}(K) \lesssim \mathcal{R}_2^{\text{KeRNS}}(K) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3.$$

where  $\mathcal{R}_2^{\text{KeRNS}}(K)$  is the UCBVI-type regret bound given in Theorem 4 for KeRNS.

*Proof.* We use the same regret decomposition as in the proof of Theorem 5, but the term (C) is bounded differently (as in Lemma 12).

Using Lemma 18, Lemma 22, and the same arguments as in the proof of Lemma 12, we have

$$\begin{aligned}
 (\mathbf{C}) &= \left[ \tilde{P}_h^k - P_h^k \right] \left( \tilde{V}_{h+1}^k - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &= \left[ \tilde{P}_h^k - \hat{P}_h^k \right] \left( \tilde{V}_{h+1}^k - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) + \left[ \hat{P}_h^k - P_h^k \right] \left( \tilde{V}_{h+1}^k - V_{k,h+1}^* \right) (\tilde{x}_h^k, \tilde{a}_h^k) \\
 &\lesssim L(\varepsilon + \varepsilon_{\mathcal{X}})H + \frac{\varepsilon}{\sigma}H^2 + \left[ \hat{P}_h^k - P_h^k \right] (V_{h+1}^k - V_{k,h+1}^*) (\tilde{x}_h^k, \tilde{a}_h^k) \quad (\text{by lemmas 18 and 22}) \\
 &\lesssim L(\varepsilon + \varepsilon_{\mathcal{X}})H + \frac{\varepsilon}{\sigma}H^2 + \frac{1}{H}P_h^k (V_{h+1}^k - V_{k,h+1}^*) (x_h^k, a_h^k) + \frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\
 &+ L\sigma + \sum_{h'=h}^H \mathbf{bias}(k, h') + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \quad (\text{following the proof of Lemma 12}) \\
 &\lesssim L(\varepsilon + \varepsilon_{\mathcal{X}})H + \frac{\varepsilon}{\sigma}H^2 + \frac{1}{H}P_h^k \left( \tilde{V}_{h+1}^k - V_{k,h+1}^* \right) (x_h^k, a_h^k) + \frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \\
 &+ L\sigma + \sum_{h'=h}^H \mathbf{bias}(k, h') + L\rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)]
 \end{aligned}$$

where, in the last line, we used the fact that  $V_{k,h+1}^* \leq V_{k,h+1}^*$  and that  $V_{k+1}^k \leq \tilde{V}_{k+1}^k + L(\varepsilon + \varepsilon_{\mathcal{X}})H + (\varepsilon/\sigma)H^2$  by Lemma 22.

Putting together the bounds for (A) – (D) and using the same arguments as in the proof of Theorem 5, especially the inequalities (8), (9) and (10), we obtain

$$\begin{aligned}
 \sum_{k=1}^K \check{\delta}_1^k &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} + H^2 |\mathcal{C}_{\sigma}| \\
 &+ \sum_{k=1}^K \sum_{h=1}^H \left( 1 + \frac{1}{H} \right)^{H-h+1} \tilde{\xi}_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH\sigma + L(\varepsilon + \varepsilon_{\mathcal{X}})H^2 + \frac{\varepsilon}{\sigma}H^3
 \end{aligned}$$

where  $\tilde{\xi}_{h+1}^k$  is a martingale difference sequence with respect to  $(\mathcal{F}_h^k)_{k,h}$  bounded by  $4H$ .

Now, from lemmas 10 and 22, we have:

$$\begin{aligned}
 \mathcal{R}^{\text{RS-KeRNS}}(K) &= \sum_{k=1}^K \left( V_{k,1}^*(x_1^k) - V_{k,h}^*(x_1^k) \right) \leq \sum_{k=1}^K \left( V_1^k(x_1^k) - V_{k,h}^*(x_1^k) \right) + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) \\
 &\lesssim \sum_{k=1}^K \left( \tilde{V}_1^k(x_1^k) - V_{k,h}^*(x_1^k) \right) + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma}KH^2 \\
 &= \sum_{k=1}^K \check{\delta}_1^k + \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH(\varepsilon + \varepsilon_{\mathcal{X}}) + \frac{\varepsilon}{\sigma}KH^2 \\
 &\lesssim \sum_{k=1}^K \sum_{h=1}^H \left( \frac{H}{\sqrt{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)}} + \frac{H^2 |\mathcal{C}'_{\sigma}|}{\mathbf{C}_h^k(\tilde{x}_h^k, \tilde{a}_h^k)} \right) \mathbb{I} \{ \rho [(x_h^k, a_h^k), (\tilde{x}_h^k, \tilde{a}_h^k)] \leq 2\sigma \} + H^2 |\mathcal{C}_{\sigma}| \\
 &+ \sum_{k=1}^K \sum_{h=1}^H \left( 1 + \frac{1}{H} \right)^h \tilde{\xi}_{h+1}^k + H \sum_{k=1}^K \sum_{h=1}^H \mathbf{bias}(k, h) + LKH\sigma + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3
 \end{aligned}$$

Then, as in the proof of Theorem 4, we obtain

$$\mathcal{R}^{\text{RS-KeRNS}}(K) \lesssim \mathcal{R}_2^{\text{KeRNS}}(K) + L(\varepsilon + \varepsilon_{\mathcal{X}})KH^2 + \frac{\varepsilon}{\sigma}KH^3$$



with probability at least  $1 - \delta$ .

□

## H Technical Lemmas

**Lemma 23** (adapted from Domingues et al. (2020)). Consider a sequence of non-negative real numbers  $\{z_s\}_{s=1}^t$  and let  $\bar{\Gamma}_{(\eta, W)} : \mathbb{R}_+ \rightarrow [0, 1]$  satisfy Assumption 4. For a given  $t$ , let

$$w_s \stackrel{\text{def}}{=} \bar{\Gamma}_{(\eta, W)}\left(t - s - 1, \frac{z_s}{\sigma}\right) \quad \text{and} \quad \tilde{w}_s \stackrel{\text{def}}{=} \frac{w_s}{\beta + \sum_{s'=1}^t w_{s'}}$$

for  $\beta > 0$ . Then, we have

$$\sum_{s=1}^t \tilde{w}_s z_s \leq 2\sigma \left(1 + \sqrt{\log(C_1 t / \beta + e)}\right).$$

*Proof.* For completeness, we reproduce here the proof of Lemma 7 of Domingues et al. (2020), which also applies to our setting. We split the sum into two terms:

$$\sum_{s=1}^t \tilde{w}_s z_s = \sum_{s: z_s < c} \tilde{w}_s z_s + \sum_{s: z_s \geq c} \tilde{w}_s z_s \leq c + \sum_{s: z_s \geq c} \tilde{w}_s.$$

From Assumption 4, we have  $w_s \leq C_1 \exp(-z_s^2/(2\sigma^2))$ . Hence,  $\tilde{w}_s \leq (C_1/\beta) \exp(-z_s^2/(2\sigma^2))$ , since  $\beta + \sum_{s'=1}^t w_{s'} \geq \beta$ .

We want to find  $c$  such that:

$$z_s \geq c \implies \frac{C_1}{\beta} \exp\left(-\frac{z_s^2}{2\sigma^2}\right) \leq \frac{1}{t} \frac{2\sigma^2}{z_s^2}$$

which implies, for  $z_s \geq c$ , that  $\tilde{w}_s \leq \frac{1}{t} \frac{2\sigma^2}{z_s^2}$ .

Let  $x = z_s^2/2\sigma^2$ . Reformulating, we want to find a value  $c'$  such that  $C_1 \exp(-x) \leq \beta/(xt)$  for all  $x \geq c'$ . Let  $c' = 2 \log(C_1 t / \beta + e)$ . If  $x \geq c'$ , we have:

$$\begin{aligned} \frac{x}{2} \geq \log\left(\frac{C_1 t}{\beta} + e\right) &\implies x \geq \frac{x}{2} + \log\left(\frac{C_1 t}{\beta} + e\right) \implies x \geq \log x + \log(C_1 t / \beta + e) \\ &\implies (C_1 / \beta) \exp(-x) \leq 1/(xt) \end{aligned}$$

as we wanted.

Now,  $x \geq c'$  is equivalent to  $z_s \geq \sqrt{2\sigma^2 c'} = 2\sigma \sqrt{\log(C_1 t / \beta + e)}$ . Therefore, we take  $c = 2\sigma \sqrt{\log(C_1 t / \beta + e)}$ , which gives us

$$\sum_{s: z_s \geq c} \tilde{w}_s z_s \leq \sum_{s: z_s \geq c} \frac{1}{t} \frac{2\sigma^2}{z_s^2} z_s \leq \frac{2\sigma^2}{t} \sum_{s: z_s \geq c} \frac{1}{z_s} \leq \frac{2\sigma^2}{c} \frac{|\{s : z_s \geq c\}|}{t} \leq \frac{2\sigma^2}{c}.$$

Finally, we obtain:

$$\begin{aligned} \sum_{s=1}^t \tilde{w}_s z_s &\leq c + \sum_{s: z_s \geq c} \tilde{w}_s z_s \leq c + \frac{2\sigma^2}{c} \\ &= 2\sigma \sqrt{\log(C_1 t / \beta + e)} + \frac{\sigma}{\sqrt{\log(C_1 t / \beta + e)}} \leq 2\sigma \left(1 + \sqrt{\log(C_1 t / \beta)}\right). \end{aligned}$$

□

**Lemma 24.** Let  $\bar{\Gamma}_{(\eta, W)} : \mathbb{R}_+ \rightarrow [0, 1]$  be a kernel that satisfies Assumption 4. Let  $a \in \mathbb{R}_+^t$  and  $f_1, f_2, f_3$  be functions from  $\mathbb{R}_+^t$  to  $\mathbb{R}$  defined as

$$\begin{aligned} f_1(z) &= \frac{\sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma) a_s}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)}, \\ f_2(z) &= \sqrt{\frac{1}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)}}, \\ f_3(z) &= \frac{1}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)}. \end{aligned}$$

Then, for any  $y, z \in \mathbb{R}_+$ , we have

$$\begin{aligned} |f_1(z) - f_1(y)| &\leq \frac{2C_2 \|a\|_\infty t}{\beta\sigma} \|z - y\|_\infty \\ |f_2(z) - f_2(y)| &\leq \frac{C_2 t}{2\beta^{3/2}\sigma} \|z - y\|_\infty \\ |f_3(z) - f_3(y)| &\leq \frac{C_2 t}{\beta^2\sigma} \|z - y\|_\infty \end{aligned}$$

*Proof.* From Assumption 4, the function  $z \mapsto \bar{\Gamma}_{(\eta, W)}(t-s-1, z)$  is  $C_2$ -Lipschitz, which yields

$$\begin{aligned} &|f_1(z) - f_1(y)| \\ &\leq \left| \frac{\sum_{s=1}^t (\bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma) - \bar{\Gamma}_{(\eta, W)}(t-s-1, y_s/\sigma)) a_s}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)} \right| \\ &+ \left| \frac{\sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, y_s/\sigma) a_s}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, y_s/\sigma)} \right| \left| \frac{\sum_{s=1}^t (\bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma) - \bar{\Gamma}_{(\eta, W)}(t-s-1, y_s/\sigma))}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)} \right| \\ &\leq \frac{C_2 \sum_{s=1}^t (1/\sigma) |z_s - y_s| a_s}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)} + \|a\|_\infty \frac{C_2 \sum_{s=1}^t (1/\sigma) |z_s - y_s|}{\beta + \sum_{s=1}^t \bar{\Gamma}_{(\eta, W)}(t-s-1, z_s/\sigma)} \\ &\leq \frac{2C_2 \|a\|_\infty t}{\beta\sigma} \|z - y\|_\infty. \end{aligned}$$

The proofs for  $f_2$  and  $f_3$  are analogous. For  $f_2$ , we also use the fact that the function  $x \mapsto (1/\sqrt{\beta+x})$  is  $1/(2\beta^{3/2})$ -Lipschitz.  $\square$

**Lemma 25.** Let  $a \in \mathbb{R}_+^t$  and  $f_1, f_2, f_3$  be functions from  $\mathbb{R}_+^t$  to  $\mathbb{R}$  defined as

$$\begin{aligned} f_1(z) &= \frac{\sum_{s=1}^t g(s) \exp(-z_s^2/(2\sigma^2)) a_s}{\beta + \sum_{s=1}^t g(s) \exp(-z_s^2/(2\sigma^2))}, \\ f_2(z) &= \sqrt{\frac{1}{\beta + \sum_{s=1}^t g(s) \exp(-z_s^2/(2\sigma^2))}}, \\ f_3(z) &= \frac{1}{\beta + \sum_{s=1}^t g(s) \exp(-z_s^2/(2\sigma^2))} \end{aligned}$$

where  $g : \mathbb{N}^* \mapsto [0, 1]$  is an arbitrary function bounded by 1. Then, for any  $y, z \in \mathbb{R}_+$ , we have

$$\begin{aligned} |f_1(z) - f_1(y)| &\leq \frac{4 \|a\|_\infty}{\sigma} \left(1 + \sqrt{\log^+(t/\beta)}\right) \|z - y\|_\infty \\ |f_2(z) - f_2(y)| &\leq \frac{1}{2\beta^{1/2}\sigma} \left(1 + \sqrt{\log^+(t/\beta)}\right) \|z - y\|_\infty \\ |f_3(z) - f_3(y)| &\leq \frac{1}{\beta\sigma} \left(1 + \sqrt{\log^+(t/\beta)}\right) \|z - y\|_\infty \end{aligned}$$

*Proof.* We use the fact that, for any differentiable  $f : \mathbb{R}_+^t \rightarrow \mathbb{R}$ ,

$$|f(x) - f(y)| \leq \sup_{z \in \mathbb{R}_+^t} \|\nabla f(z)\|_1 \|x - y\|_\infty.$$

We have

$$\left| \frac{\partial f_1(z)}{\partial z_i} \right| \leq \frac{2 \|a\|_\infty}{\sigma^2} \frac{g(i) z_i \exp(-z_i^2 / (2\sigma^2))}{\beta + \sum_{s=1}^t g(s) z_s \exp(-z_s^2 / (2\sigma^2))}$$

which implies, by Lemma 23,

$$\begin{aligned} \|\nabla f_1(z)\|_1 &\leq \frac{2 \|a\|_\infty}{\sigma^2} \frac{\sum_{i=1}^t g(i) z_i \exp(-z_i^2 / (2\sigma^2))}{\beta + \sum_{s=1}^t g(s) z_s \exp(-z_s^2 / (2\sigma^2))} \\ &\leq \frac{2 \|a\|_\infty}{\sigma^2} 2\sigma \left( 1 + \sqrt{\log^+(t/\beta)} \right) = \frac{4 \|a\|_\infty}{\sigma} \left( 1 + \sqrt{\log^+(t/\beta)} \right). \end{aligned}$$

The proofs for  $f_2$  and  $f_3$  are analogous.  $\square$

**Lemma 26** (value functions are Lipschitz continuous). *Under Assumptions 1 and 2, for all  $(k, h)$ , the functions  $V_{k,h}^*$  and  $Q_{k,h}^*$  are  $L_h$ -Lipschitz, where  $L_h \stackrel{\text{def}}{=} \sum_{h'=h}^H L_r L_p^{H-h'}$ .*

*Proof.* This fact is proved in Lemma 4 of Domingues et al. (2020) and also in Proposition 2.5 of Sinclair et al. (2019). For completeness, we also present a proof here.

We proceed by induction. For  $h = H$ ,  $Q_{k,H}^*(x, a) = r_H^k(x, a)$  which is  $L_r$ -Lipschitz by Assumption 2. Also,

$$V_{k,H}^*(x) - V_{k,H}^*(y) = \max_a Q_{k,H}^*(x, a) - \max_a Q_{k,H}^*(y, a) \leq \max_a (Q_{k,H}^*(x, a) - Q_{k,H}^*(y, a)) \quad (11)$$

$$\leq \max_a L_H \rho[(x, a), (y, a)] \leq L_H \rho_{\mathcal{X}}(x, y), \quad \text{by Assumption 1} \quad (12)$$

which verifies the induction hypothesis for  $h = H$ , since we can invert the roles of  $x$  and  $y$  to obtain  $|V_{k,H}^*(x) - V_{k,H}^*(y)| \leq L_H \rho_{\mathcal{X}}(x, y)$ .

Now, assume that the hypothesis is true for  $h + 1$ , i.e., that  $V_{k,h+1}^*$  and  $Q_{k,h+1}^*$  are  $L_{h+1}$ -Lipschitz. We have

$$\begin{aligned} Q_{k,h}^*(x, a) - Q_{k,h}^*(x', a') &\leq L_r \rho[(x, a), (x', a')] + \int_{\mathcal{X}} V_{k,h+1}^*(y) (P_h(dy|x, a) - P_h(dy|x', a')) \\ &\leq L_r \rho[(x, a), (x', a')] + L_{h+1} \int_{\mathcal{X}} \frac{V_{k,h+1}^*(y)}{L_{h+1}} (P_h(dy|x, a) - P_h(dy|x', a')) \\ &\leq \left[ L_r + L_p \sum_{h'=h+1}^H L_r L_p^{H-h'} \right] \rho[(x, a), (x', a')] \\ &= \sum_{h'=h}^H L_r L_p^{H-h'} \rho[(x, a), (x', a')] \end{aligned}$$

where, in last inequality, we use fact that  $V_{k,h+1}^*/L_{h+1}$  is 1-Lipschitz, the definition of the 1-Wasserstein distance and Assumption 2. The same argument used in Eq. 11 shows that  $|V_{k,h}^*(x) - V_{k,h}^*(y)| \leq L_h \rho_{\mathcal{X}}(x, y)$ , which concludes the proof.  $\square$

**Fact 1** (small useful result).

$$|\min(a, b) - \min(a, c)| \leq |b - c|$$

*Proof.* Since  $\min(x, y) = (x + y)/2 - |x - y|/2$ , we have

$$\begin{aligned}\min(a, b) - \min(a, c) &= \frac{a + b}{2} - \frac{|a - b|}{2} - \frac{a + c}{2} + \frac{|a - c|}{2} \\ &= \frac{b - c}{2} + \frac{|a - c| - |a - b|}{2} \\ &\leq \frac{b - c}{2} + \frac{|b - c|}{2} \\ &\leq |b - c|\end{aligned}$$

where we used the fact that  $|a - c| \leq |a - b| + |b - c|$ . By symmetry,  $\min(a, c) - \min(a, b) \leq |b - c|$ . which gives us the result.  $\square$

## I Experiments

### I.1 Setup

We consider a continuous MDP whose state-space is the unit ball in  $\mathbb{R}^2$  with four actions, representing a move to the right, left, up or down. Each action results in a displacement of 0.1 in the corresponding direction, plus a Gaussian noise, in both directions, of standard variation 0.01. The agent starts at  $(0, 0)$ . Let  $b_i^k \in \{0, 0.25, 0.5, 0.75, 1\}$  and  $x_i \in \{(0.8, 0.0), (0.0, 0.8), (-0.8, 0.0), (0.0, -0.8)\}$ . We consider the following mean reward function:

$$r_h^k(x, a) = \sum_{i=1}^4 b_i^k \max\left(0, 1 - \frac{\|x - x_i\|_2}{0.5}\right)$$

Every  $N$  episodes, the coefficients  $b_i^k$  are changed, which impact the optimal policy.

Taking  $\eta = \exp(-(1/N)^{2/3})$ , we tested the Gaussian kernel  $\Gamma(t, u, v) = \eta^t \exp(-\rho[u, v]^2 / (2\sigma^2))$  and a higher-order kernel  $\Gamma(t, u, v) = \eta^t \exp(-(\rho[u, v] / \sigma)^4 / 2)$ . We set  $\sigma = 0.05$ ,  $\varepsilon = \varepsilon_{\mathcal{X}} = 0.1$ ,  $\beta = 0.01$ ,  $H = 15$ .

We ran the experiment with horizon  $H = 15$  for  $2 \times 10^4$  episodes. Every  $N$  episodes, the coefficients  $b_i^k$  were changed, according to Table 4.

Table 4: Value of  $b_i^k$  according to  $x_i$  and the episode  $k$

episode / $x_i$	(0.8, 0.0)	(0.0, 0.8)	(-0.8, 0.0)	(0.0, -0.8)
$k \bmod N = 0$	1/4	0	0	0
$k \bmod N = 1$	1/4	1/2	0	0
$k \bmod N = 2$	1/4	1/2	3/4	0
$k \bmod N = 3$	1/4	1/2	3/4	1

We took  $\beta = 0.01$  and used the following simplified exploration bonuses:

$$B_h^k(x, a) = \frac{0.1}{\sqrt{C_h^k(x, a)}} + \frac{\beta H}{C_h^k(x, a)} \quad (13)$$

where the factor 0.1 was chosen in order to ensure that the baseline is able to learn a good policy in less than 1000 episodes, i.e., before there is a change in the environment.

Additionally, to take into account the fact that the Lipschitz constant is rarely known in practical problems, we replaced the interpolation step (line 8 of Alg. 5) by a nearest-neighbor search in the representative states:

$$\tilde{Q}_h^k(x, a) = \tilde{Q}_{h, \zeta}^k(x', a'), \quad \text{where } (x', a') = \underset{(\bar{x}, \bar{a}) \in \bar{\mathcal{X}}_h^k \times \bar{\mathcal{A}}_h^k}{\operatorname{argmin}} \rho[(x, a), (\bar{x}, \bar{a})].$$

### I.2 Results

Figures 2 and 3 show the total reward and the regret of **RS-KeRNS** compared to baselines for the two choices of kernel function (Gaussian and 4-th order kernel), for 3 different values of  $\Delta$ , which is determined by the period  $N$  of changes in the MDP (the reward changes every  $N$  episodes).

In all experiments we observe that **Kernel-UCBVI** is not able to adapt to the changes in the environment, whereas **RS-KeRNS** is able to track the behavior of the baseline **RestartBaseline** which knows when the changes happen and resets the reward estimator when there is a change.

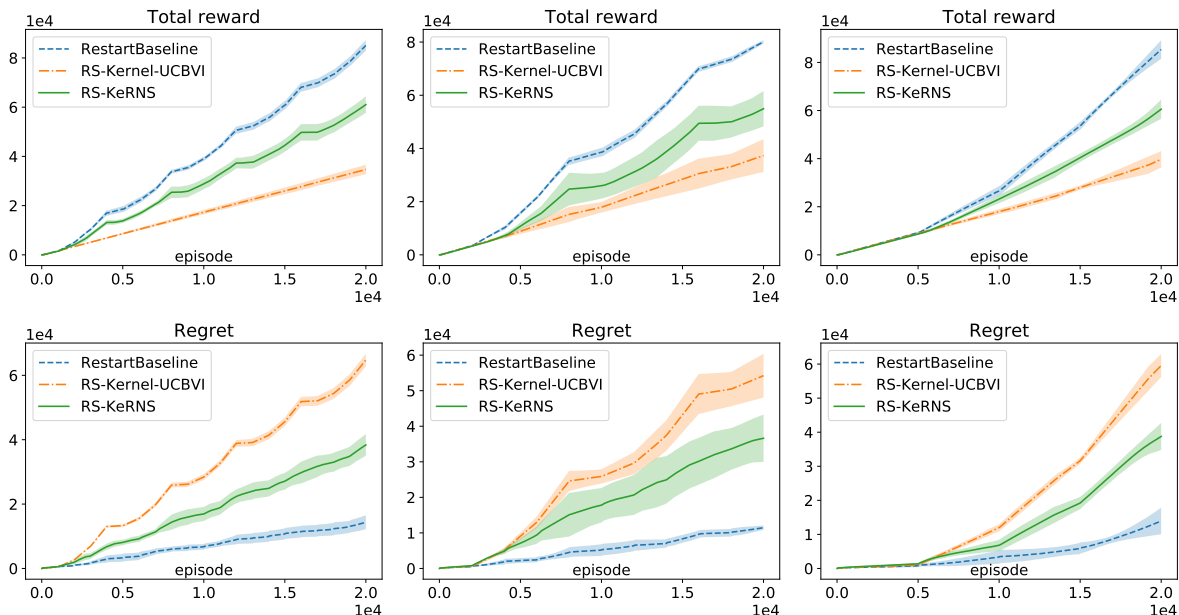


Figure 2: Total reward (top row) and regret(bottom row) of **RS-KeRNS** compared to baselines, using the Gaussian kernel  $\Gamma(t, u, v) = \eta^t \exp(-\rho[u, v]^2 / (2\sigma^2))$ . The figures on the left, in the middle, and on the right correspond to  $N = 1000$ ,  $N = 2000$  and  $N = 5000$ , respectively, where  $N$  is the period of the changes in the MDP. Average over 4 runs.

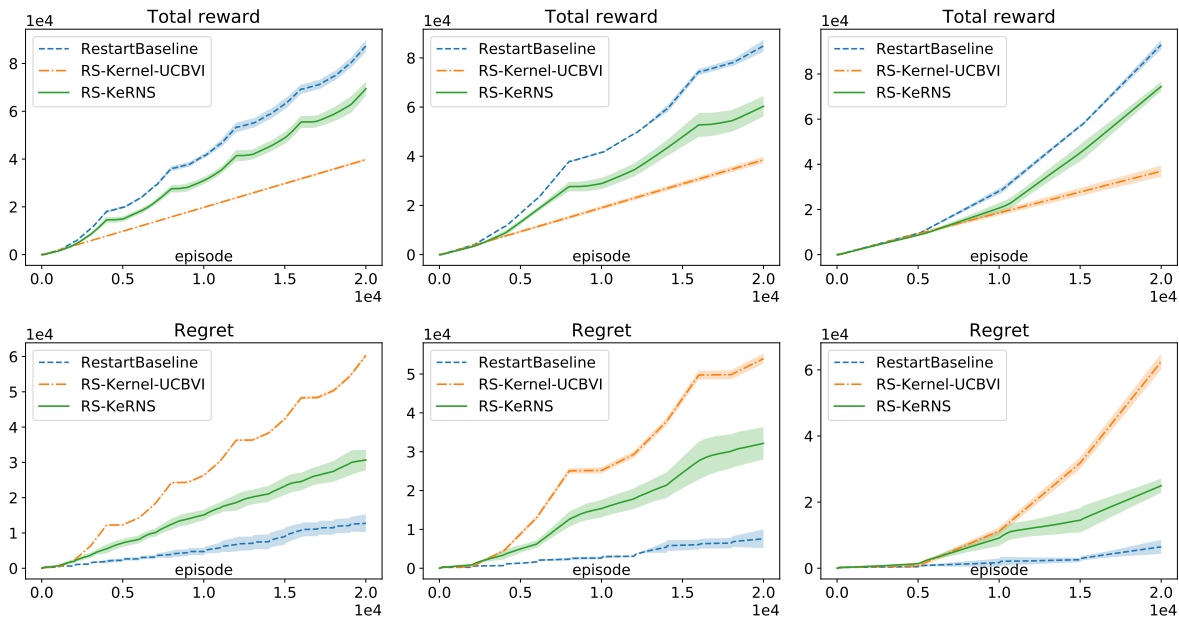


Figure 3: Total reward (top row) and regret(bottom row) of **RS-KeRNS** compared to baselines, using the kernel  $\Gamma(t, u, v) = \eta^t \exp(-(\rho[u, v] / \sigma)^4 / 2)$ . The figures on the left, in the middle, and on the right correspond to  $N = 1000$ ,  $N = 2000$  and  $N = 5000$ , respectively, where  $N$  is the period of the changes in the MDP. Average over 4 runs.