
Model-Based Reinforcement Learning with Value-Targeted Regression

Alex Ayoub¹ Zeyu Jia² Csaba Szepesvári^{1,3} Mengdi Wang^{4,3} Lin F. Yang⁵

Abstract

This paper studies model-based reinforcement learning (RL) for regret minimization. We focus on finite-horizon episodic RL where the transition model P belongs to a known family of models \mathcal{P} , a special case of which is when models in \mathcal{P} take the form of linear mixtures: $P_\theta = \sum_{i=1}^d \theta_i P_i$. We propose a model based RL algorithm that is based on the optimism principle: In each episode, the set of models that are ‘consistent’ with the data collected is constructed. The criterion of consistency is based on the total squared error that the model incurs on the task of predicting *state values* as determined by the last value estimate along the transitions. The next value function is then chosen by solving the optimistic planning problem with the constructed set of models. We derive a bound on the regret, which, in the special case of linear mixtures, takes the form $\tilde{O}(d\sqrt{H^3T})$, where H , T and d are the horizon, the total number of steps and the dimension of θ , respectively. In particular, this regret bound is independent of the total number of states or actions, and is close to a lower bound $\Omega(\sqrt{HdT})$. For a general model family \mathcal{P} , the regret bound is derived based on the Eluder dimension.

1. Introduction

In reinforcement learning (RL), a core problem in artificial intelligence (Russel & Norvig, 2003; Sutton & Barto, 2018), an agent learns to control a possibly complex, initially unknown environment in a sequential trial and error process. The application of RL algorithms to various do-

mains, such as games, robotics and science, has witnessed phenomenal empirical advances during the last few years (e.g., Mnih et al. 2015; Silver et al. 2017; AlQuraishi 2019; Arulkumaran et al. 2019). In online RL, an agent has to learn to act in an unknown environment “from scratch”, collect data as she acts, and adapt her policy to maximize the reward collected, or, equivalently, to minimize her regret. Designing RL algorithms that provably achieve sublinear regret in some class of environments has been the subject of much research, mainly focusing on the so-called tabular, and linear-factored MDP settings (e.g., Jaksch et al. 2010; Osband et al. 2014; Azar et al. 2017; Dann et al. 2017; 2018; Agrawal & Jia 2017; Osband et al. 2017; Jin et al. 2018; Yang & Wang 2019a; Jin et al. 2019). An appealing alternative to studying these structured cases is to consider learning and acting when the environment is described by a general model class, the topic of the current paper. Despite its appeal, as it appears, prior work, considered this option exclusively in a Bayesian setting. In particular, (Strens, 2000) introduced posterior sampling to RL, which was later analyzed by (Osband & Van Roy, 2014; Abbasi-Yadkori & Szepesvári, 2015; Theocharous et al., 2017) (for a more in-depth discussion of related work, the reader is referred to Section 5). As opposed to these works, in the present paper we are interested in developing algorithms for bounding the worst-case expected regret.

The specific setting that we adopt is that of episodic reinforcement learning in an environment where the unknown transition probability model that describes the environment’s stochastic dynamics belongs to a family of models that is given to the learner. The model family \mathcal{P} is a general set of models, and it may be either finitely parametrized or nonparametric. In particular, our approach accommodates working with smoothly parameterized models (e.g., Abbasi-Yadkori & Szepesvári, 2015), and can find use in both robotics (Kober et al., 2013) and queuing systems (Kovalenko, 1968). An illuminating special case is when elements of \mathcal{P} take the form $P_\theta = \sum_i \theta_i P_i$ where P_1, P_2, \dots, P_d are fixed, known basis models and $\theta = (\theta_1, \dots, \theta_d)$ are unknown, real-valued parameters. Model P_θ can be viewed as a linear mixture model that aggregates a finite family of known basic dynamical models (Modi et al., 2019). As an important special case, linear mixture models include the linear-factor MDP model of

¹Amii and Department of Computing Science, University of Alberta ²School of Mathematical Science, Peking University ³DeepMind ⁴Department of Electrical Engineering, Princeton University ⁵Department of Electrical and Computer Engineering, University of California, Los Angeles. Correspondence to: Csaba Szepesvári <csaba.szepesvari@gmail.com>, Mengdi Wang <mengdiw@princeton.edu>.

(Yang & Wang, 2019a).

The main contribution of this paper is a new model-based upper confidence RL algorithm. The main novelty is the criterion to select models that are deemed consistent with past data. As opposed to the most common approach where the models are selected based on their ability to predict next states or raw observations (cf. Jaksch et al. 2010; Yang & Wang 2019a or Strens 2000; Osband & Van Roy 2014; Abbasi-Yadkori & Szepesvári 2015; Ouyang et al. 2017; Agrawal & Jia 2017 in a Bayesian setting), we propose to evaluate models based on their ability to predict the values of a value function at next states, where the value function used is an estimate of the optimal that our algorithm produces based on past information. In essence, our algorithm selects models based on their ability to produce small prediction errors in an appropriately constructed *value-targeted* regression (VTR) problem. Our algorithm combines VTR for constructing sets of plausible models with (standard) optimistic planning. The idea of using a value function estimate to “fit” models has been explored in the context of batch RL by Farahmand (2018).

VTR is attractive for multiple reasons: (i) Firstly, VTR permits model learning to *focus on task-relevant aspects* of the transition dynamics. This is important as the dynamics can be quite complicated and in a resource bounded setting, modelling irrelevant aspects of the dynamics can draw valuable resources away from modelling task-relevant aspects. (ii) Secondly, VTR poses model learning as a real-valued regression problem, which should be easier than the usual approaches to build probabilistic models. In particular, when state-representation available to the agent takes values in a high-dimensional space then building a faithful probability model can be highly demanding. (iii) Thirdly, VTR aims to control directly what matters in terms of controlling the regret. Specifically the objective used in value-targeted is obtained from an expression that upper bounds the regret, hence it is natural to expect that minimizing value prediction errors will lead to a small regret.

An additional attractive feature of our algorithm is its modular structure. As a result, advances on the components (faster optimistic planning, tighter confidence sets for VTR) are directly translated into an improved algorithm. On the skeptical side, one may question whether VTR is going “too far” in ignoring details of the dynamics. In particular, since the value functions used in defining the regression targets are derived based on imperfect knowledge, the model may never get sufficiently refined in a way that would keep the regret small. Similarly, one may be worried about that by ignoring details of the observations (i.e., the identity of states), the approach advocated is ignoring information available in the data, which may slow down learning. This leads to central question of our article:

Is value-targeted regression sufficient and efficient for model-based online RL?

The main contribution of this paper is a positive answer to this question. In particular, the regret bounds we derive conclusively show that despite the imperfection and non-stationarity of the value targets, our algorithm will not get “stuck” (i.e., it enjoys sublinear regret). Our results further suggest that in the worst case sense, for common settings, there may be no performance penalty associated with using value-targeted regression. We are careful here as this conclusion is based on comparing worst-case upper bounds, which cannot provide a definitive answer. Finally, it is worth noting that the regret bound does not depend on the size of either the state *or* the action space.

To complement the theoretical findings, results from a number of small-scale, synthetic experiments confirm that our algorithm is competitive in terms of its regret. The experiments also allow us to conclude that it is value-targeted regression *together* with optimistic planning that is effective. In particular, if optimism is taken away (i.e., ϵ -greedy for exploration), value-targeted regression performs worse than using a canonical approach to estimate the model. Similarly, if value-targeted regression is taken away, optimism together with the canonical model-estimation approach is less effective. Note that our results do not rule out that certain combinations of value-targeted regression and canonical model building are more effective than value-targeted regression. In fact, given the vast number of possibilities, we find this to be quite probable. While our proofs can be adjusted to deal with adding simultaneous alternative targets, sadly, our current theoretical tools are unable to capture the tradeoffs that one expects to arise as a result of such modifications.

It is interesting to note that, in an independent and concurrent work, value-targeted regression has also been suggested as the main model building tool of the MuZero algorithm (Schrittwieser et al., 2019), which was empirically evaluated on a number of RL benchmarks, such as the 57 Atari “games”, the game of “Go”, chess and shogi, and was found to be highly competitive with its state-of-the-art alternatives. This reinforces the conclusion that training models using value-targeted regression is indeed a good approach to build effective model-based RL algorithms. Since MuZero does *not* implement optimistic planning and our results show that optimism is not optional in a worst case sense, the good results of MuZero on these benchmark may seem to contradict our experimental findings that value-targeted regression is ineffective without an appropriate, ‘smart’ exploration component. However, there is no contradiction: Smart exploration may be optional in some environments; our experiments show that it is not optional on some environments. In short, for robust performance across a wide range of environments, smart exploration is necessary but

smart exploration may be optional in some environments.

2. Problem Formulation

We study learning to control episodic Markov decision processes (MDPs, for short), described by a tuple $M = (\mathcal{S}, \mathcal{A}, P, r, H, s_\circ)$. Here, \mathcal{S} is the state space, \mathcal{A} is the action space, P is the transition kernel, r is a reward function, $H > 0$ is the episode length, or horizon, and $s_\circ \in \mathcal{S}$ is the initial state. In the online RL problem, the learning agent is given \mathcal{S} , \mathcal{A} , H and r but does not know P .¹ The agent interacts with its environment in a number of episodes. Each episode begins at state s_\circ and ends after the agent made H decisions. At state $s \in \mathcal{S}$, the agent, after observing the state s , can choose an action $a \in \mathcal{A}$. As a result, the immediate reward $r(s, a)$ is incurred. Then the process transitions to a random next state $s' \in \mathcal{S}$ according to the transition law $P(\cdot|s, a)$.² The agent's goal is to maximize the total expected reward received over time.

If P is known, the behavior that achieves maximum expected reward over any number of episodes can be described by applying a deterministic policy π . Such a policy is a mapping from $\mathcal{S} \times [H]$ into \mathcal{A} (note for a natural number n , $[n] = \{1, \dots, n\}$). Following the policy means that the agent upon encountering state s in stage h will choose action $\pi(s, h)$. In what follows, we will use $\pi_h(s)$ as an alternate notation, as this makes some of the formulae more readable. (We will follow the same convention of moving h to the subindex position when it comes to other functions whose domain is $\mathcal{S} \times [H]$.)

The value function $V^\pi : \mathcal{S} \times [H] \rightarrow \mathbb{R}$ of a policy π is defined via

$$V_h^\pi(s) = \mathbb{E}_\pi \left[\sum_{i=h}^H r(s_i, \pi(s_i)) \mid s_h = s \right], \quad s \in \mathcal{S},$$

where s_i is the state encountered at stage $i \in [H]$ and the subscript π (which we will often suppress) signifies that the probabilities underlying the expectation are jointly governed by π and P (P is suppressed for clarity). An optimal policy π^* and the optimal value function V^* are defined to be a policy and its value function such that $V_h^\pi(s)$ with $\pi = \pi^*$ achieves the maximum value among all possible policies for any $s \in \mathcal{S}$ and $h \in [H]$. Note that this is well-defined and in fact, as noted above, there is no loss of generality in restricting the search of optimal policies to deterministic policies.

¹Our results are easy to extend to the case when r is not known.

²The precise definitions require measure-theoretic concepts (Bertsekas & Shreve, 1978), i.e., P is a Markov kernel, mapping from $\mathcal{S} \times \mathcal{A}$ to distributions over \mathcal{S} , hence, all these spaces need to be properly equipped with a measurability structure. For the sake of readability and also because they are well understood, we omit these technical details.

In online RL, a learning agent will use all past observations to come up with its decisions. The performance of such an agent is measured by its regret, which is the total reward the agent misses because she did not follow the optimal policy from the beginning. In particular, the total expected regret of an agent \mathcal{A} across K episodes is given by

$$R(T) = \mathbb{E} \left[\sum_{k=1}^K \left(V_1^*(s_1^k) - \sum_{h=1}^H r(s_h^k, a_h^k) \right) \right], \quad (1)$$

where $T = KH$ is the total number of time steps that the agent interacts with its environment, $s_1^k = s_\circ$ is the initial state at the start of the k -th episode, and $s_1^1, a_1^1, \dots, s_H^k, a_H^k, s_1^2, a_1^2, \dots, s_H^2, a_H^2, s_1^K, a_1^K, \dots, s_H^K, a_H^K$ are the $T = KH$ state-action pairs in the order that they are encountered by the agent. The regret is sublinear if $R(T)/T \rightarrow 0$ as $T \rightarrow \infty$. As is well known, the worst-case value of $R(T)$ over a set of sufficiently large model class, grows at least as fast as \sqrt{T} regardless of the algorithm used (e.g., Jaksch et al., 2010).

In this paper, we aim to design a general model-based reinforcement learning algorithm, with a guaranteed sublinear regret, for any (not too large) family of transition models:

Assumption 1 (Known Transition Model Family). *The unknown transition model P belongs to a family of models \mathcal{P} which is available to the learning agent. The elements of \mathcal{P} are transition kernels mapping state-action pairs to signed distributions over \mathcal{S} .*

That we allow signed distributions only increases generality; this may be important when one is given a model class that can be compactly represented but only when it also includes non-probability kernels (see Pires & Szepesvári 2016 for a discussion of this). Parametric and nonparametric transition models are common in modelling complex stochastic controlled systems. For example, robotic systems are often smoothly parameterized by unknown mechanical parameters such as friction and parameters that describe the geometry of the robot.

An important special case is the class of linear mixture models:

Definition 1 (Linear Mixture Models). *We say that \mathcal{P} is the class of linear mixture models with component models P_1, \dots, P_d if P_1, \dots, P_d are transition kernels that map state-action pairs to signed measures and $P \in \mathcal{P}$ if and only if there exists $\theta \in \mathbb{R}^d$ such that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,*

$$P(ds'|s, a) = \sum_{j=1}^d \theta_j P_j(ds'|s, a) = P.(ds'|s, a)^\top \theta. \quad (2)$$

The linear mixture model can be viewed as a way of aggregating a number of known basis models as considered by

(Modi et al., 2019). We can view each $P_j(\cdot|\cdot)$ as a basis latent “mode”. When θ is restricted to lie in the $(d - 1)$ simplex, the actual transition is a probabilistic mixture of these latent modes. As an example of when mixture models arise, consider large-scale queueing networks where the arrival rate and job processing speed for each queue is not known. By using a discrete-time Bernoulli approximation, the transition probability matrix from time t to $t + \Delta t$ becomes increasingly close to linear with respect to the unknown arrival/processing rates as $\Delta t \rightarrow 0$. In this case, it is common to model the discrete-time state transition as a linear aggregation of arrival/processing processes with unknown parameters (Kovalenko, 1968).

Another interesting special case is the linear-factored MDP model of (Yang & Wang, 2019a) where, assuming a discrete state space for a moment, P takes the form

$$\begin{aligned} P(s'|s, a) &= \phi(s, a)^\top M \psi(s') \\ &= \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} M_{ij} [\psi_j(s') \phi_i(s, a)], \end{aligned}$$

where $\phi(s, a) \in \mathbb{R}^{d_1}$, $\psi(s') \in \mathbb{R}^{d_2}$ are given features for every $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$ (when the state space is continuous, ψ becomes an \mathbb{R}^{d_2} -valued measure over \mathcal{S}). The matrix $M \in \mathbb{R}^{d_1 \times d_2}$ is an unknown matrix and is to be learned. It is easy to see that the factored MDP model is a special case of the linear mixture model (19) with each $\psi_j(s') \phi_i(s, a)$ being a basis model (this should be replaced by $\psi_j(ds') \phi_i(s, a)$ when the state space is continuous). In this case, the number of unknown parameters in the transition model is $d = d_1 \times d_2$. In this setting, without additional assumptions, our regret bound will match that of (Yang & Wang, 2019a).

3. Upper Confidence RL with Value-Targeted Model Regression

Our algorithm (Alg 1) can be viewed as a generalization of UCRL (Jaksch et al., 2010), following ideas of (Osband & Van Roy, 2014). In particular, at the beginning of episode $k = 1, 2, \dots, K$, the algorithm first computes a subset B_k of the model class \mathcal{P} that contains the set of models that are deemed to be consistent with all the data that has been collected in the past. The new idea, value-targeted regression is used in the construction of B_k . The details of how this is done are postponed to a later section. Given B_k , the algorithm needs to find the model that maximizes the optimal value, and the corresponding optimal policy. Denoting by V_P^* the optimal value function under a model P , this amounts to finding the model $P \in B_k$ that maximizes the value $V_{P,1}^*(s_1^k)$. Given the model P_k that maximizes this value, an optimal policy is extracted from the model as in standard dynamic programming, detailed in the next section.

Algorithm 1 UCRL-VTR

- 1: **Input:** Family of MDP models \mathcal{P} , $d, H, T = KH$, sequence $\{\beta_k\}_{k=1,2,\dots}$
- 2: $B_1 = \mathcal{P}$
- 3: **for** $k = 1, 2, \dots, K$ **do**
- 4: Observe the initial state s_1^k of episode k
- 5: **Optimistic planning:**

$$P_k = \operatorname{argmax}_{P' \in B_k} V_{P',1}^*(s_1^k)$$

 Compute $Q_{1,k}, \dots, Q_{H,k}$ for P_k using (3)

- 6: **for** $h = 1, 2, \dots, H$ **do**
- 7: Choose the next action greedily with respect to $Q_{h,k}$:

$$a_h^k = \operatorname{arg max}_{a \in \mathcal{A}} Q_{h,k}(s_h^k, a)$$
- 8: Observe state s_{h+1}^k
- 9: Compute and store value predictions: $y_{h,k} \leftarrow V_{h+1,k}(s_{h+1}^k)$
- 10: **end for**
- 11: **Construct confidence set using value-targeted regression as described in Section 3.2:**

$$B_{k+1} = \{P' \in \mathcal{P} | L_{k+1}(P', \hat{P}_{k+1}) \leq \beta_k\}$$

- 12: **end for**
-

At the end of the episode, the data collected is used to refine the confidence set B_k .

3.1. Model-Based Optimistic Planning

Upper confidence methods are prominent in sequential online learning. As noted before, we let

$$P_k = \operatorname{argmax}_{P' \in B_k} V_{P',1}^*(s_1^k).$$

Given model P_k , the optimal policy for P_k can be computed using dynamic programming. In particular, for $1 \leq h \leq H + 1$, define

$$\begin{aligned} Q_{H+1,k}(s, a) &= 0, \\ V_{h,k}(s) &= \max_{a \in \mathcal{A}} Q_{h,k}(s, a), \\ Q_{h,k}(s, a) &= r(s, a) + \langle P_k(\cdot|s, a), V_{h+1,k} \rangle, \end{aligned} \tag{3}$$

where for a measure μ and function f that share a common domain, $\langle \mu, f \rangle$ denotes the integral of f with respect to μ . It follows that, taking the action at state s and stage h that maximizes $Q_{h,k}(s, \cdot)$ gives an optimal policy for model P_k . As long as $P \in B_k$ with high probability, the preceding calculation gives an optimistic (that is, upper) estimate of value of an episode. Next, we show how to construct the confidence set B_k .

3.2. Value-Targeted Regression for Confidence Set Construction

Every time we observe a transition (s, a, s') with $s' \sim P(\cdot|s, a)$, we receive information about the model P . A standard approach to use this information would be either using a maximum likelihood approach, or regressing “onto” s' . As our goal is not to find the best model, we propose an alternate approach where we set up a regression problem where the model is used to predict the value assigned to s' by our more recent value function estimate:

$$\hat{P}_{k+1} = \operatorname{argmin}_{P' \in \mathcal{P}} \sum_{k'=1}^k \sum_{h=1}^H L'(P'), \quad (4)$$

$$L'(P') = \left(\langle P'(\cdot|s_h^{k'}, a_h^{k'}), V_{h+1, k'} \rangle - y_{h, k'} \right)^2$$

$$y_{h, k'} = V_{h+1, k'}(s_{h+1}^{k'}), \quad h \in [H], k' \in [k].$$

In the above regression procedure, the regret target keeps changing as the algorithm refines the value estimates. This is in contrast to typical supervised learning for building models, where the regression targets are often fixed objects (such as raw observations, features or keypoints; e.g. (Jaksch et al., 2010; Osband & Van Roy, 2014; Abbasi-Yadkori & Szepesvári, 2015; Xie et al., 2016; Agrawal & Jia, 2017; Yang & Wang, 2019a; Kaiser et al., 2019)). For a confidence set construction, we get inspiration from Proposition 5 in the paper of (Osband & Van Roy, 2014). The set is centered at \hat{P}_{k+1} . Defining

$$L_{k+1}(P, \hat{P}_{k+1}) = \sum_{k'=1}^k \sum_{h=1}^H \left(\langle P(\cdot|s_h^{k'}, a_h^{k'}), \hat{P}_{k+1}(\cdot|s_h^{k'}, a_h^{k'}) \rangle - V_{h+1, k'} \right)^2$$

we let

$$B_{k+1} = \{P' \in \mathcal{P} \mid L_{k+1}(P', \hat{P}_{k+1}) \leq \beta_{k+1}\},$$

where the value of β_k is obtained using a calculation similar to that done in Proposition 5 of the paper of (Osband & Van Roy, 2014). In turn, this calculation is based on the nonlinear least-squares confidence set construction from Russo & Van Roy (2014), which we describe and refine in the appendix. It is not hard to see that the confidence set can also be written in the alternative form

$$B_{k+1} = \{P' \in \mathcal{P} \mid \tilde{L}_{k+1}(P') \leq \tilde{\beta}_{k+1}\}$$

with a suitably defined $\tilde{\beta}_{k+1}$ and where

$$\tilde{L}_{k+1}(P') = \sum_{k'=1}^k \sum_{h=1}^H \left(\langle P'(\cdot|s_h^{k'}, a_h^{k'}), V_{h+1, k'} \rangle - y_{h, k'} \right)^2.$$

Note that the above formulation strongly exploits that the MDP is time-homogeneous: The same transition model is used at all stages of an episode. When the MDP is time-inhomogeneous, the construction can be easily modified to accommodate this.

3.3. Implementation of UCRL-VTR

Algorithm 1 gives a general and modular template for model-based RL that is compatible with regression methods/optimistic planners. While the algorithm is conceptually simple, and the optimization and evaluation of the loss in value-targeted regression appears to be at advantage in terms of computation as compared to standard approaches typically used in model-based RL, the implementation of UCRL-VTR is nontrivial in general and for now it requires a case-by-base design.

Computation efficiency of the algorithm depends on the specific family of models chosen. For the linear-factor MDP model considered by Yang & Wang (2019a), the regression is linear and admits efficient implementation; further, optimistic planning for this model can be implemented in $\text{poly}(d)$ time by using Monte-Carlo simulation and sketching as argued in the cited paper. Other ideas include loosening the confidence set to come up with computationally tractable methods, or relaxing the requirement that the same model is used in all stages. This latter idea is what we use in our experiments. In the general case, optimistic planning is computationally intractable. However, we expect that randomized (e.g. (Osband et al., 2017; 2014; Lu & Van Roy, 2017)) and approximate dynamic programming methods (tree search, roll out, see e.g., (Bertsekas & Tsitsiklis, 1996)) will often lead to tractable and good approximations. As was mentioned above, in some special cases these have been rigorously shown to work. In similar settings, the approximation errors are known to mildly impact the regret (Abbasi-Yadkori & Szepesvári, 2015) and we expect the same to hold in our setting. If we look beyond methods with rigorous guarantees, there are practical deep RL algorithms that implement parts of UCRL-VTR. As mentioned earlier, the MuZero algorithm of Schrittwieser et al. (2019) is a state-of-the-art algorithm on multiple domains and this algorithm implements value-targeted-regression to learn a model which is fed to a planner that uses Monte Carlo tree search, although the planner does not implement optimistic planning.

4. Theoretical Analysis

We will need the concept of Eluder dimension due to (Russo & Van Roy, 2014). Let \mathcal{F} be a set of real-valued functions with domain \mathcal{X} . For $f \in \mathcal{F}$, $x_1, \dots, x_t \in \mathcal{X}$, introduce the notation $f|_{(x_1, \dots, x_t)} = (f(x_1), \dots, f(x_t))$. We say that $x \in \mathcal{X}$ is ϵ -independent of $x_1, \dots, x_t \in \mathcal{X}$ given \mathcal{F} if there exists $f, f' \in \mathcal{F}$ such that $\|(f - f')|_{(x_1, \dots, x_t)}\|_2 \leq \epsilon$ while $f(x) - f'(x) > \epsilon$.

Definition 2 (Eluder dimension, Russo & Van Roy 2014). The Eluder dimension $\dim_{\mathcal{E}}(\mathcal{F}, \epsilon)$ of \mathcal{F} at scale ϵ is the length of the longest sequence (x_1, \dots, x_n) in \mathcal{X} such that for some $\epsilon' \geq \epsilon$, for any $2 \leq t \leq n$, x_t is ϵ' -independent of

(x_1, \dots, x_{t-1}) given \mathcal{F} .

Let \mathcal{V} be the set of optimal value functions under some model in \mathcal{P} : $\mathcal{V} = \{V_{P'}^* : P' \in \mathcal{P}\}$. Note that $\mathcal{V} \subset \mathcal{B}(\mathcal{S}, H)$, where $\mathcal{B}(\mathcal{S}, H)$ denotes the set of real-valued measurable functions with domain \mathcal{S} that are bounded by H . We let $\mathcal{X} = \mathcal{S} \times \mathcal{A} \times \mathcal{V}$. Choose \mathcal{F} to be the collection of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ as follows:

$$\mathcal{F} = \left\{ f \mid \begin{array}{l} \exists P \in \mathcal{P} \text{ s.t. for any } (s, a, v) \in \mathcal{S} \times \mathcal{A} \times \mathcal{V} \\ f(s, a, v) = \int P(ds' | s, a) v(s') \end{array} \right\}. \quad (5)$$

Note that $\mathcal{F} \subset \mathcal{B}(\mathcal{X}, H)$. For a norm $\|\cdot\|$ on \mathcal{F} and $\alpha > 0$ let $\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|)$ denote the $(\alpha, \|\cdot\|)$ -covering number of \mathcal{F} . That is, this if $m = \mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|)$ then one can find m elements of \mathcal{F} such that any element in \mathcal{F} is at most α away from one of these elements in norm $\|\cdot\|$. Denote by $\|\cdot\|_\infty$ the supremum norm: $\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$.

Define the K -episode *pseudo-regret* as

$$R_K = \sum_{k=1}^K (V^*(s_0^k) - V^{\pi_k}(s_0^k)).$$

Clearly, $R(KH) = \mathbb{E}R_K$ holds for any $K > 0$ where $R(T)$ is the expected regret after T steps of interaction as defined in (1). Thus, to study the expected regret, it suffices to study R_K . Our main result is as follows.

Theorem 1 (Regret of Algorithm 1). *Let Assumption 1 hold and let $\alpha \in (0, 1)$. For $k > 0$ let β_k be*

$$\begin{aligned} \beta_k &= 2H^2 \log \left(\frac{2\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) \\ &+ 2H(kH - 1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{4kH(kH - 1)}{\delta} \right)} \right\}. \end{aligned} \quad (6)$$

Then, for any $K > 0$, with probability $1 - 2\delta$,

$$\begin{aligned} R_K &\leq \alpha + H(d \wedge K(H - 1)) + 4\sqrt{d\beta_K}K(H - 1) \\ &+ H\sqrt{2K(H - 1)\log(1/\delta)}, \end{aligned}$$

where $d = \dim_{\mathcal{E}}(\mathcal{F}, \alpha)$ is the Eluder dimension with \mathcal{F} given by (5).

A typical choice of α is $\alpha = 1/(KH)$. In the special case of linear transition model, Theorem 1 implies a worst-case regret bound that depends linearly on the number of parameters.

Corollary 2 (Regret of Algorithm 1 for Linearly-Parametrized Transition Model). *Let P_1, \dots, P_d be d transition models, $\Theta \subset \mathbb{R}^d$ a nonempty set with diameter R measured in $\|\cdot\|_1$ and let $\mathcal{P} = \{\sum_j \theta_j P_j : \theta \in \Theta\}$. Then, for any $0 < \delta < 1$, with probability at least $1 - \delta$, the*

pseudo-regret R_K of Algorithm 1 when it uses the confidence sets given in Theorem 1 satisfies

$$R_K = \tilde{O}(d\sqrt{H^3 K \log(1/\delta)}).$$

We also provide a lower bound for the regret in our model. The proof is by reduction to a known lower bound and is left to Appendix B.

Theorem 3 (Regret Lower Bound). *For any $H \geq 1$ and $d \geq 8$, there exist a state space \mathcal{S} and action set \mathcal{A} , a reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, d transition models P_1, \dots, P_d and a set $\Theta \subset \mathbb{R}^d$ of diameter of at most one such that for any algorithm there exists $\theta \in \Theta$ such that for sufficiently large number of episodes K , the expected regret of the algorithm on the H -horizon MDP with reward r and transition model $P = \sum_j \theta_j P_j$ is at least $\Omega(H\sqrt{dK})$.*

Rusmevichientong & Tsitsiklis (2010) gave a regret lower bound of $\Omega(d\sqrt{T})$ for linearly parameterized bandit with actions on the unit sphere (see also Section 24.2 of Lattimore & Szepesvári 2020). Our regret upper bound matches this bandit lower bound in d, T . Whether the upper or lower bound is tight (or none of them) remains to be seen. The theorems validate that, in the setting we consider, it is sufficient to use the predicted value functions as regression targets. That for the special case of linear mixture models the lower bound is close to the upper bound appears to suggest that little benefit if any can be derived from fitting the transition model to predict future observations. We conjecture that this is in fact true when considering the worst-case regret. Of course, a conclusion that is concerned with the worst-case regret has no implication for the behavior of the respective methods on particular MDP instances. We note in passing that by appropriately increasing β_k , the regret upper bounds can be extended to the so-called misspecified case when P can be outside of \mathcal{P} (for related results, see, e.g., Jin et al. 2019; Lattimore & Szepesvári 2019). However, the details of this are left for future work.

Further, our method applies to handle the case where the linearly parameterized transition model is sparse. Suppose that model parameter θ is known to have at most s nonzero entries. In this case, the class of sparse linear models has a much smaller covering number, and the regret would improve and depend on both d, s . Details of this are left for future work.

4.1. Regret Bound with Model Misspecification

Next, we consider the case where the model family \mathcal{P} does not exactly realize the true transition model P :

Assumption 2 (Model with misspecification error). *The model family \mathcal{P} ε -approximates P in the sense that there exists $P^* \in \mathcal{P}$ such that*

$$\sup_{s,a} \|P(\cdot|s, a) - P^*(\cdot|s, a)\|_{TV} \leq \varepsilon, \quad (7)$$

where $\|\cdot\|_{TV}$ denotes the total variation distance.

This assumption indicates that the true transition model P of the MDP is close to the family \mathcal{P} , and ε measures the worst-case deviation. We handle the misspecification error by slightly enlarging the confidence set. This allows us to obtain a regret bound similar to our previous result with an additional linear term that is proportional to the misspecification error and slightly larger constants:

Theorem 4. *Let Assumption 2 hold, $\alpha, \delta \in (0, 1)$. We choose β_k be*

$$\begin{aligned} \beta_k &= 8H^2 \log \left(\frac{4\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) \\ &+ 4H(kH - 1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{8kH(kH - 1)}{\delta} \right)} \right\} \\ &+ 8H^3 k\varepsilon^2. \end{aligned}$$

Then, for any $K > 0$, with probability $1 - 2\delta$, the K -episode pseudo-regret R_K of Algorithm 1 satisfies

$$\begin{aligned} R_K &\leq \alpha + H(d \wedge K(H - 1)) + 4\sqrt{d\beta_K K(H - 1)} \\ &+ H\sqrt{2K(H - 1)\log(1/\delta)} + H^2 K\varepsilon, \end{aligned}$$

where $d = \dim_{\mathcal{E}}(\mathcal{F}, \alpha)$ with \mathcal{F} given by (5).

The proof of this theorem is given in Appendix D.

5. Related Work

A number of prior efforts have established efficient RL methods with provable regret bounds. For tabular H -horizon MDP with S states and A actions, there have been results on model-based methods (e.g., Jaksch et al. 2010; Osband et al. 2014; Azar et al. 2017; Dann et al. 2017; Agrawal & Jia 2017; Dann et al. 2018; Kakade et al. 2018), and on model-free methods (e.g., Osband et al. 2017; Jin et al. 2018; Zhang et al. 2020). Both model-based and model-free methods are known to achieve a regret of $\tilde{O}(\sqrt{H^2 SAT})$, where $\tilde{O}(\cdot)$ hides log factors, T denotes the total number of timesteps and the bound applies to the non-stationary setting (when the transition models are not shared across stages). Moreover, apart from logarithmic factors, this bound is known to be unimprovable Jaksch et al. (2010); Osband et al. (2017); Jin et al. (2018); Zhang et al. (2020). There have been significant theoretical and empirical advances on RL with function approximation, including but not limited to (Baird, 1995; Tsitsiklis & Van Roy, 1997; Parr et al., 2008; Mnih et al., 2013; 2015; Silver et al., 2017; Yang & Wang, 2019b; Bradtke & Barto, 1996). Under the assumption that the optimal action-value function is captured by linear features, Zanette et al. (2019) considers the case when the features are ‘‘extrapolation friendly’’ and a simulation

oracle is available, Wen & Van Roy (2013; 2017) tackle problems where the transition model is deterministic, Du et al. (2019) deals with a relaxation of the deterministic case when the transition model has low variance. Yang & Wang (2019b) considers the case of linear factor models, while Lattimore & Szepesvari (2019) considers the case when all the action-value functions of all deterministic policies are well-approximated using a linear function approximator. These latter works handle problems when the algorithm has access to a simulation oracle of the MDP. As for regret minimization in RL using linear function approximation, Yang & Wang (2019a) assumed the transition model admits a matrix embedding of the form $P(s'|s, a) = \phi(s, a)^\top M\psi(s')$, and proposed a model-based MatrixRL method with regret bounds $\tilde{O}(H^2 d\sqrt{T})$ with stronger assumptions and $\tilde{O}(H^2 d^2\sqrt{T})$ in general, where d is the dimension of state representation $\phi(s, a)$. Jin et al. (2019) studied the setting of linear-factor MDP and constructed a model-free least-squares action-value iteration algorithm, which was proved to achieve the regret bound $\tilde{O}(\sqrt{H^3 d^3 T})$. Modi et al. (2019) considered a related setting where the transition model is an ensemble involving state-action-dependent features and basis models and proved a sample complexity $\frac{d^3 K^2 H^2}{\varepsilon^2}$ where d is the feature dimension, K is the number of basis models and $d \cdot K$ is their total model complexity. As for RL with a general model class, in their seminal work, Osband & Van Roy (2014) provided a general posterior sampling RL method that works for any given classes of reward and transition functions. It established a Bayesian regret upper bound $O(\sqrt{d_K d_E T})$, where d_K and d_E are the Kolmogorov and the Eluder dimensions of the model class. In the case of linearly parametrized transition model (Assumption 1 of this paper), this Bayesian regret becomes $O(d\sqrt{T})$, and our worst-case regret result matches with the Bayesian one. Abbasi-Yadkori & Szepesvari (2015); Theodorou et al. (2017) also considered the Bayesian regret and, in particular, Abbasi-Yadkori & Szepesvari (2015) considered a smooth parameterization. To the authors’ best knowledge, there are no prior works addressing the problem of designing low-regret algorithms for MDPs with linearly or non-linearly parameterized transition models. Model based PAC RL algorithms have been studied by Sun et al. (2019), who essentially adopt the value-aware loss of Farahmand et al. (2017), who considered this loss in the batch setting. Farahmand (2018) refines the work of Farahmand et al. (2017) by changing the algorithm to be similar to what we use here: In every iteration of fitted Q-iteration, first a model is obtained by minimizing the value prediction loss measured with the last value function, after which this model is used to obtain the next action-value function. The main result bounds the suboptimality of the policy that is greedy with respect to the last action-value function. A preliminary version of our paper was presented at L4DC.

Exploration/ Targets	Optimism	Dithering
Next states	UC-MatrixRL	EG-Freq
Values	UCRL-VTR	EG-VTR
Mixed	UCRL-Mixed	EG-Mixed

Table 1. Legend to the algorithms compared. Note that UC-MatrixRL of Yang & Wang (2019a) in the tabular case essentially becomes UCRL of Jaksch et al. (2010).

6. Numerical Experiments

The goal of our experiments is to provide insight into the benefits and/or pitfalls of using value-targets for fitting models, both with and without optimistic planning. We run our experiments in the tabular setting as it is easy to keep aspects of the test environments under control while avoiding approximate computations. Note that tabular environments are a special case of the linear model where $P_j(s'|s, a) = \mathbb{I}(j = f(s, a, s'))$, where $j \in [S^2A]$ and f is a bijection that maps its arguments to the set $[S^2A]$, making $d = S^2A$. The objective is either to minimize mean-squared error of predicting next states (alternatively, maximize log-likelihood of observed data), which leads to frequency based model estimates, or it is to minimize the value targets as proposed in our paper. The other component of the algorithms is whether they implement optimistic planning, or planning with the nominal model and then implementing an ϵ -greedy policy with respect to the estimated model (“dithering”). In the case of optimistic planning, the algorithm that uses mixed targets uses a union bound and takes the smallest value upper confidence bounds amongst the two bounds obtained with the two model-estimation methods. This leads to six algorithms, as shown in Table 1. Results for the “mixed” variants are very similar to the variant that uses VTR and can be found in Appendix G. In the experiments we use confidence bounds that are specialized to the linear case. For details of these, see Appendix F. For ϵ -greedy, we optimize the value of ϵ in each environment to get the best results. This gives ϵ -greedy an “unfair advantage”. As we shall see it soon, despite this advantage, ϵ -greedy will not fair particularly well in our experiments.

6.1. Measurements

We report the cumulative regret as a function of the number of episodes and the weighted model error to indicate how well the model is learned. The results are obtained from 30 independent runs for the ϵ -greedy algorithms and 10 independent runs for the UC algorithms. The weighted model error reported is as follows. Given the model estimate

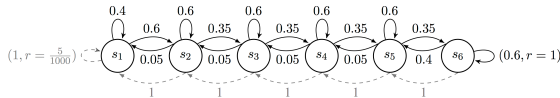


Figure 1. The “RiverSwim” environment with 6 states. State s_1 has a small associated reward, state s_6 has a large associated reward. The action whose effect is shown with the dashed arrow deterministically “moves the agent” left. The other action is stochastic, and with relatively high probability moves the agent towards state s_6 : This represents swimming “against the current”.

\hat{P} , we compute

$$E(\hat{P}) = \sum_{s,a} \sum_{s'} \frac{N(s, a, s')}{N'(s, a)} |\hat{P}(s' | s, a) - P^*(s' | s, a)|, \quad (8)$$

where N' is the observation-count of the state-action pair (s, a) , N is the count of transitioning to s' from (s, a) , and P^* is the true dynamics model. The weighting is introduced so that an algorithm that discards a state-action pair is not (unduly) penalized.

6.2. Results for RiverSwim

The schematic diagram of the RiverSwim environment is shown in Figure 1. RiverSwim consists of S states arranged in a chain. The agent begins on the far left and has the choice of swimming left or right at each state. There is a current that makes swimming left much easier than swimming right. Swimming left with the current always succeeds in moving the agent left, but swimming right against the current sometimes moves the agent right (with a small probability of moving left as well), but more often than not leaves the agent in the current state. Thus smart exploration is a necessity to learn a good policy in this environment. We experiment with small environments with $S = 5$ and set the horizon to $H = 20$. The optimal value of the initial state is 5.6 for our five-state RiverSwim. The initial state is the leftmost state (s_1 in the diagram). The value that we found to work the best for EGRL-VTR is $\epsilon = 0.2$ and the value that we found to work best for EG-Freq is $\epsilon = 0.12$. The results are shown in Figure 2. The regret per episode for an algorithm that “does not learn” is expected to be in the same range as the respective optimal values. Based on this we see that 10^5 episodes is barely sufficient for the algorithms other than UCRL-VTR to learn a good policy. Looking at the model errors we see that EGRL-VTR is doing quite poorly, EG-Freq is also lacking, the others are doing reasonably well. However, this is because EG-Freq visits more uniformly than the other methods the various state-action pairs. The results clearly indicate that (i) fitting to the state-value function alone provides enough of a signal for learning as evident by UCRL-VTR obtaining low regret as predicted by our theoretical results, and that

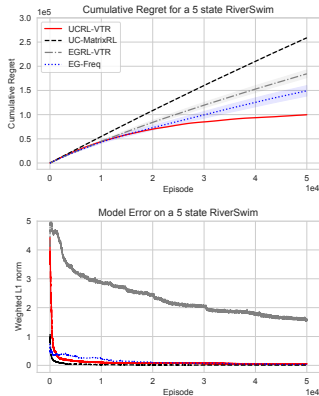


Figure 2. The results for the ϵ -greedy algorithms are averaged over thirty runs and error bars are reported for the regret plots.

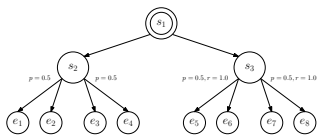


Figure 3. An eleven state WideTree MDP. The algorithm starts in the initial state s_1 . From the initial state s_1 the algorithm has a choice of either deterministically transitioning to either state s_2 or state s_3 . Finally from either state s_2 or state s_3 the algorithm picks one of two possible actions and transitions to one of the terminal states e_i . The choice of the initial action determines the delayed reward the algorithm will observe.

(ii) optimism is necessary when using VTR to achieve good results, as evident by UCRL-VTR achieving significantly better regret than EGRL-VTR and even in the smaller RiverSwim environment. It is also promising that value-targeted regression with optimistic exploration outperformed optimism based on the “canonical” model estimation procedure. We attribute this to the fact that value-targeted regression will learn a model faster that predicts the optimal values well than the canonical, frequency based approach. That value-targeted regression also learns a model with small weighted error appears to be an accidental feature of this environment. Our next experiments are targeted at further exploring whether VTR can be effective *without* learning a good model.

6.3. WideTree

We introduce a novel tabular MDP we call WideTree. The WideTree environment has a fixed horizon $H = 2$ and $S = 11$ states. A visualization of the WideTree environment is shown in Figure 3. In WideTree, an agent starts at the initial state s_1 . The agent then progresses to one of the many bottom terminal states and collects a reward of either 0 or 1. The only significant action is whether to transition from s_1 to either s_2 or s_3 . Note that the model in the second layer is

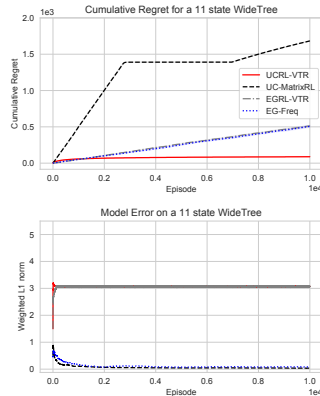


Figure 4. The results for the ϵ -greedy algorithms are averaged over thirty runs and error bars are reported for the regret plots.

irrelevant for making a good decision: Once in s_3 , all actions lead to a reward of one, and once in s_2 , all actions lead to a reward of zero. We vary the number of bottom states reachable from states s_2 and s_3 while still maintaining a reward structure but the results here are shown with $S = 11$. We set $\epsilon = 0.1$ in this environment, as this allows the model error of EG-Freq to match that of UC-MatrixRL. The results are shown in Figure 4. Both UCRL-VTR and EG-VTR learn equally poor models (their graphs are ‘on the top of each other’). Yet, UCRL-VTR manages to quickly learn a good policy, as attested by its low regret. EG-Freq and EG-VTR perform equally poorly and UC-MatrixRL is even slower as it keeps exploring the environment. These experiments clearly illustrate that UCRL-VTR is able to achieve good results without learning a good model – its focus on values makes pays off swiftly in this well-chosen environment.

7. Conclusions

We considered online learning in episodic MDPs and proposed an optimistic model-based reinforcement learning method (UCRL-VTR) with the unique characteristic to evaluate and select models based on their ability to predict value functions that the algorithm constructs during learning. The regret of the algorithm was shown to be bounded by a quantity that relates to the richness of the model class through the Eluder dimension and the metric entropy of an appropriately constructed function space. For the case of linear mixture models, the regret bound simplifies to $\tilde{O}(d\sqrt{H^3T})$ where d is the number of model parameters, H is the horizon, and T is the total number of interaction steps. Our experiments confirmed that the value-targeted regression objective is not only theoretically sound, but also yields a competitive method which allows task-focused model-tuning: In a carefully chosen environment we demonstrated that the algorithm achieves low regret despite that it ignores modeling a major part of the environment.

8. Acknowledgements

Csaba Szepesvári gratefully acknowledges funding from the Canada CIFAR AI Chairs Program, Amii and NSERC. Mengdi Wang gratefully acknowledges funding from the U.S. National Science Foundation (NSF) grant CMMI-1653435, Air Force Office of Scientific Research (AFOSR) grant FA9550-19-1-020, and C3.ai DTI.

References

- Abbasi-Yadkori, Y. and Szepesvári, C. Bayesian optimal control of smoothly parameterized systems. In *UAI*, pp. 1–11, 2015.
- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Agrawal, S. and Jia, R. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In *Advances in Neural Information Processing Systems*, pp. 1184–1194, 2017.
- AlQuraishi, M. AlphaFold at CASP13. *Bioinformatics*, 35(22):4862–4865, 2019.
- Arulkumaran, K., Cully, A., and Togelius, J. Alphastar: An evolutionary computation perspective. *arXiv preprint arXiv:1902.01724*, 2019.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 263–272. JMLR. org, 2017.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pp. 30–37. Elsevier, 1995.
- Bertsekas, D. P. and Shreve, S. *Stochastic optimal control: the discrete-time case*. Academic Press, 1978.
- Bertsekas, D. P. and Tsitsiklis, J. N. *Neuro-dynamic programming*. Athena Scientific, 1996.
- Bradtke, S. J. and Barto, A. G. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- Dann, C., Lattimore, T., and Brunskill, E. Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5713–5723, 2017.
- Dann, C., Li, L., Wei, W., and Brunskill, E. Policy certificates: Towards accountable reinforcement learning. *arXiv preprint arXiv:1811.03056*, 2018.
- Du, S. S., Luo, Y., Wang, R., and Zhang, H. Provably efficient Q -learning with function approximation via distribution shift error checking oracle. *arXiv preprint arXiv:1906.06321*, 2019.
- Farahmand, A.-M. Iterative value-aware model learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 9090–9101, 2018.
- Farahmand, A.-M., Barreto, A., and Nikovski, D. Value-aware loss function for model-based reinforcement learning. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1486–1494. PMLR, 2017.
- Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is Q -learning provably efficient? In *Advances in Neural Information Processing Systems*, pp. 4863–4873, 2018.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019.
- Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Koza-kowski, P., Levine, S., Mohiuddin, A., Sepassi, R., Tucker, G., and Michalewski, H. Model-based reinforcement learning for Atari. In *ICLR*, 2019.
- Kakade, S., Wang, M., and Yang, L. F. Variance reduction methods for sublinear reinforcement learning. *arXiv preprint arXiv:1802.09184*, 2018.
- Kober, J., Bagnell, J. A., and Peters, J. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Kovalenko, B. G. I. N. *Introduction to queueing theory*. Israel Program for Scientific Translation, Jerusalem, 1968.
- Lattimore, T. and Szepesvári, C. Learning with good feature representations in bandits and in RL with a generative model, 2019.
- Lattimore, T. and Szepesvári, C. *Bandit Algorithms*. Cambridge University Press, 2020. (to appear).
- Lu, X. and Van Roy, B. Ensemble sampling. In *Advances in neural information processing systems*, pp. 3258–3266, 2017.

- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529, 2015.
- Modi, A., Jiang, N., Tewari, A., and Singh, S. Sample complexity of reinforcement learning using linearly combined model ensembles. *arXiv preprint arXiv:1910.10597*, 2019.
- Osband, I. and Van Roy, B. Model-based reinforcement learning and the Eluder dimension. In *Advances in Neural Information Processing Systems*, pp. 1466–1474, 2014.
- Osband, I. and Van Roy, B. On lower bounds for regret in reinforcement learning. *arXiv preprint arXiv:1608.02732*, 2016.
- Osband, I., Van Roy, B., and Wen, Z. Generalization and exploration via randomized value functions. *arXiv preprint arXiv:1402.0635*, 2014.
- Osband, I., Van Roy, B., Russo, D., and Wen, Z. Deep exploration via randomized value functions. *arXiv preprint arXiv:1703.07608*, 2017.
- Ouyang, Y., Gagrani, M., Nayyar, A., and Jain, R. Learning unknown Markov decision processes: A Thompson sampling approach. In *Advances in Neural Information Processing Systems*, pp. 1333–1342, 2017.
- Parr, R., Li, L., Taylor, G., Painter-Wakefield, C., and Littman, M. L. An analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pp. 752–759. ACM, 2008.
- Pires, B. and Szepesvári, C. Policy error bounds for model-based reinforcement learning with factored linear models. In *COLT*, pp. 121–151, 2016.
- Rusmevichientong, P. and Tsitsiklis, J. N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Russel, S. and Norvig, P. *Artificial Intelligence – a modern approach*. Prentice Hall, 2003.
- Russo, D. and Van Roy, B. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., Lillicrap, T., and Silver, D. Mastering Atari, Go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*, 2019.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., Driessche, G. v. d., Graepel, T., and Hassabis, D. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- Strens, M. J. A. A Bayesian framework for reinforcement learning. In *ICML*, pp. 943–950, 2000.
- Sun, W., Jiang, N., Krishnamurthy, A., Agarwal, A., and Langford, J. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory (COLT)*, pp. 2898–2933, 2019.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. The MIT Press, 2 edition, 2018.
- Theocharous, G., Wen, Z., Abbasi-Yadkori, Y., and Vlassis, N. Posterior sampling for large scale reinforcement learning. *arXiv preprint arXiv:1711.07979*, 2017.
- Tsitsiklis, J. N. and Van Roy, B. Analysis of temporal-difference learning with function approximation. In *Advances in neural information processing systems*, pp. 1075–1081, 1997.
- Wen, Z. and Van Roy, B. Efficient exploration and value function generalization in deterministic systems. In *Advances in Neural Information Processing Systems*, pp. 3021–3029, 2013.
- Wen, Z. and Van Roy, B. Efficient reinforcement learning in deterministic systems with value function generalization. *Mathematics of Operations Research*, 42(3): 762–782, 2017.
- Xie, C., Patil, S., Moldovan, T., Levine, S., and Abbeel, P. Model-based reinforcement learning with parametrized physical models and optimism-driven exploration. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 504–511. IEEE, 2016.
- Yang, L. F. and Wang, M. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019a.
- Yang, L. F. and Wang, M. Sample-optimal parametric Q -learning with linear transition models. *International Conference on Machine Learning*, 2019b.

Zanette, A., Lazaric, A., Kochenderfer, M. J., and Brunskill, E. Limiting extrapolation in linear approximate value iteration. In *Advances in Neural Information Processing Systems 32*, pp. 5616–5625. Curran Associates, Inc., 2019.

Zhang, Z., Zhou, Y., and Ji, X. Almost optimal model-free reinforcement learning via reference-advantage decomposition. *arXiv preprint arXiv:2004.10019*, April 2020.

A. Proof of Theorem 1

In this section, we provide the regret analysis of the UCRL-VTR Algorithm (Algorithm 1). We will explain the motivation for our construction of confidence sets for general nonlinear squared estimation, and establish the regret bound for a general class of transition models, \mathcal{P} .

A.1. Preliminaries

Recall that a finite horizon MDP is $M = (\mathcal{S}, \mathcal{A}, P, r, H, s_0)$ where \mathcal{S} is the state space, \mathcal{A} is the action space, $P = (P_a)_{a \in \mathcal{A}}$ is a collection of $P_a : \mathcal{S} \rightarrow M_1(\mathcal{S})$ Markov kernels, $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the reward function, $H > 0$ is the horizon and $s_0 \in \mathcal{S}$ is the initial state. For a state $s \in \mathcal{S}$ and an action $a \in \mathcal{A}$, $P_a(s)$ gives the distribution of the next state that is obtained when action a is executed in state s . For a bounded (measurable) function $V : \mathcal{S} \rightarrow \mathbb{R}$, we will use $\langle P_a(s), V \rangle$ as the shorthand for the expected value of V at a random next state s' whose distribution is $P_a(s)$.

Given any policy π (which may or may not use the history), its value function is

$$V^\pi(s) = \mathbb{E}_{\pi, \delta_s} \left[\sum_{i=1}^H r(s_i, a_i) \right],$$

where E_{π, δ_s} is the expectation operator underlying the probability measure P_{π, δ_s} induced over sequences of state-action pairs of length H by executing policy π starting at state s in the MDP M and s_h is the state visited in stage h and action a_h is the action taken in that stage after visiting s_h . For a nonstationary Markov policy $\pi = (\pi_1, \dots, \pi_H)$, we also let

$$V_h^\pi(s) = \mathbb{E}_{\pi_{h:H}, \delta_s} \left[\sum_{i=1}^{H-h+1} r(s_i, a_i) \right]$$

be the value function of π from stage h to H . Here, $\pi_{h:H}$ denotes the policy (π_h, \dots, π_H) . The optimal value function $V^* = (V_1^*, \dots, V_H^*)$ is defined via $V_h^*(s) = \max_{\pi} V_h^\pi(s)$, $s \in \mathcal{S}$.

For simplicity assume that r is known. To indicate the dependence of V^* on the transition model P , we will write $V_P^* = (V_{P,1}^*, \dots, V_{P,H}^*)$. For convenience, we define $V_{P,H+1}^* = 0$.

Algorithm 1 is an instance of the following general model-based optimistic algorithm:

Algorithm 2 Generic Algorithm 1-Schema for finite horizon problems

- 1: **Input:** \mathcal{P} – a set of transition models, K – number of episodes, s_0 – initial state
 - 2: Set $\mathcal{B}_1 = \mathcal{P}$
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: $P^k = \operatorname{argmax}_{\tilde{P} \in \mathcal{B}_k} \{V_{\tilde{P}}^*(s_0)\}$
 - 5: $V_k = V_{P^k}^*$
 - 6: $s_1^k = s_0$
 - 7: **for** $h = 1, \dots, H$ **do**
 - 8: Choose $a_h^k = \operatorname{argmax}_{a \in \mathcal{A}} r(s_h^k, a) + \langle P_a^k(s_h^k), V_{h+1,k} \rangle$
 - 9: Observe transition to s_{h+1}^k
 - 10: **end for**
 - 11: Construct \mathcal{B}_{k+1} based on $(s_1^k, a_1^k, \dots, s_H^k, a_H^k)$
 - 12: **end for**
-

Specific instances of Algorithm 2 differ in terms of how \mathcal{B}_{k+1} is constructed. In particular, UCRL-VTR uses the construction described in Section 3.2.

Recall that $V_k = (V_{1,k}, \dots, V_{H,k}, V_{H+1,k})$ (with $V_{H+1,k} = 0$) in Algorithm 2. Let π_k be the nonstationary Markov policy chosen in episode k by Algorithm 2. Let

$$R_K = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k)$$

be the pseudo-regret of Algorithm 1 for K episodes. The following standard lemma bounds the k th term of the expression on the right-hand side.

Lemma 5. *Assuming that $P \in \mathcal{B}_k$, we have*

$$V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}_{a_h^k}(s_h^k) - P_{a_h^k}(s_h^k), V_{h,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k},$$

where

$$\xi_{h+1,k} = \langle P_{a_h^k}(s_h^k), V_{h+1,k} - V_{h+1}^{\pi_k} \rangle - (V_{h+1,k}(s_{h+1}^k) - V_{h+1}^{\pi_k}(s_{h+1}^k)).$$

Note that $(\xi_{2,1}, \xi_{3,1}, \dots, \xi_{H,1}, \xi_{2,2}, \xi_{3,2}, \dots, \xi_{H,2}, \xi_{2,3}, \dots)$ is a sequence of martingale differences.

Proof. Because $P \in \mathcal{B}_k$, $V_1^*(s_1^k) \leq V_{1,k}(s_1^k)$ by the definition of the algorithm. Hence,

$$V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq V_{1,k}(s_1^k) - V_1^{\pi_k}(s_1^k).$$

Fix $h \in [H]$. In what follows we bound $V_{h,k}(s_h^k) - V_h^{\pi_k}(s_h^k)$. By the definition of π_k , P^k and a_h^k , we have

$$\begin{aligned} V_{h,k}(s_h^k) &= r(s_h^k, a_h^k) + \langle P_{a_h^k}^k(s_h^k), V_{h+1,k} \rangle \text{ and} \\ V_h^{\pi_k}(s_h^k) &= r(s_h^k, a_h^k) + \langle P_{a_h^k}(s_h^k), V_{h+1}^{\pi_k} \rangle. \end{aligned}$$

Hence,

$$\begin{aligned} V_{h,k}(s_h^k) - V_h^{\pi_k}(s_h^k) &= \langle P_{a_h^k}^k(s_h^k), V_{h+1,k} \rangle - \langle P_{a_h^k}(s_h^k), V_{h+1}^{\pi_k} \rangle \\ &= \langle P_{a_h^k}^k(s_h^k) - P_{a_h^k}(s_h^k), V_{h+1,k} \rangle + \langle P_{a_h^k}(s_h^k), V_{h+1,k} - V_{h+1}^{\pi_k} \rangle. \end{aligned}$$

Therefore, by induction, noting that $V_{H+1,k} = 0$, we get that

$$\begin{aligned} V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) &\leq \sum_{h=1}^{H-1} \langle P_{a_h^k}^k(s_h^k) - P_{a_h^k}(s_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} \\ &\leq \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}_{a_h^k}(s_h^k) - P_{a_h^k}(s_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k}. \end{aligned}$$

□

A.2. The confidence sets for Algorithm 1

The previous lemma suggests that at the end of the k th episode, the model could be estimated using

$$\hat{P}_k = \operatorname{argmin}_{\tilde{P} \in \mathcal{P}} \sum_{k'=1}^k \sum_{h=1}^{H-1} \left(\langle \tilde{P}_{a_h^{k'}}(s_h^{k'}), V_{h+1,k'} \rangle - V_{h+1,k'}(s_{h+1}^{k'}) \right)^2 \quad (9)$$

For a confidence set construction, we get inspiration from Proposition 5 in the paper of Osband & Van Roy (2014). The set is centered at \hat{P}_k :

$$\mathcal{B}_k = \{ \tilde{P} \in \mathcal{P} : L_k(\hat{P}_k, \tilde{P}) \leq \beta_k \}, \quad (10)$$

where

$$L_k(\hat{P}, \tilde{P}) = \sum_{k'=1}^k \sum_{h=1}^{H-1} \left(\langle \tilde{P}_{a_h^{k'}}(s_h^{k'}), V_{h+1,k'} \rangle - \hat{P}_{a_h^{k'}}(s_h^{k'}), V_{h+1,k'} \right)^2.$$

Note that this is the same confidence set as described in Section 3.2. To obtain the value of β_k , we now consider the nonlinear least-squares confidence set construction from Russo & Van Roy (2014). The next section is devoted to this construction.

A.3. Confidence sets for general nonlinear least-squares

Let $(X_p, Y_p)_{p=1,2,\dots}$ be a sequence of random elements, $X_p \in \mathcal{X}$ for some measurable set \mathcal{X} and $Y_p \in \mathbb{R}$. Let \mathcal{F} be a subset of the set of real-valued measurable functions with domain \mathcal{X} . Let $\mathbb{F} = (\mathbb{F}_p)_{p=0,1,\dots}$ be a filtration such that for all $p \geq 1$, $(X_1, Y_1, \dots, X_{p-1}, Y_{p-1}, X_p)$ is \mathbb{F}_{p-1} measurable and such that there exists some function $f_* \in \mathcal{F}$ such that $\mathbb{E}[Y_p | \mathbb{F}_{p-1}] = f_*(X_p)$ holds for all $p \geq 1$. The (nonlinear) least-squares predictor given $(X_1, Y_1, \dots, X_t, Y_t)$ is $\hat{f}_t = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{p=1}^t (f(X_p) - Y_p)^2$. We say that Z is conditionally ρ -subgaussian given the σ -algebra \mathbb{F} if for all $\lambda \in \mathbb{R}$, $\log \mathbb{E}[\exp(\lambda Z) | \mathbb{F}] \leq \frac{1}{2} \lambda^2 \rho^2$. For $\alpha > 0$, let N_α be the $\|\cdot\|_\infty$ -covering number of \mathcal{F} at scale α . That is, N_α is the smallest integer for which there exist $\mathcal{G} \subset \mathcal{F}$ with N_α elements such that for any $f \in \mathcal{F}$, $\min_{g \in \mathcal{G}} \|f - g\|_\infty \leq \alpha$. For $\beta > 0$, define

$$\mathcal{F}_t(\beta) = \left\{ f \in \mathcal{F} : \sum_{p=1}^t (f(X_p) - \hat{f}_t(X_p))^2 \leq \beta \right\}.$$

We have the following theorem, the proof of which is given in Section A.6.

Theorem 6. *Let \mathbb{F} be the filtration defined above and assume that the functions in \mathcal{F} are bounded by the positive constant $C > 0$. Assume that for each $s \geq 1$, $(Y_p - f_*(X_p))_p$ is conditionally σ -subgaussian given \mathbb{F}_{p-1} . Then, for any $\alpha > 0$, with probability $1 - \delta$, for all $t \geq 1$, $f_* \in \mathcal{F}_t(\beta_t(\delta, \alpha))$, where*

$$\beta_t(\delta, \alpha) = 8\sigma^2 \log(2N_\alpha/\delta) + 4t\alpha \left(C + \sqrt{\sigma^2 \log(4t(t+1)/\delta)} \right).$$

The proof follows that of Proposition 6, Russo & Van Roy (2014), with minor improvements, which lead to a slightly better bound. In particular, with our notation, Russo & Van Roy stated their result with

$$\beta_t^{\text{RvR}}(\delta, \alpha) = 8\sigma^2 \log(2N_\alpha/\delta) + 2t\alpha \left(8C + \sqrt{8\sigma^2 \log(8t^2/\delta)} \right).$$

While $\beta_t(\delta, \alpha) \leq \beta_t^{\text{RvR}}(\delta, \alpha)$, the improvement is only in terms of smaller constants.

A.4. The choice of β_k in Algorithm 1

To use this result in our RL problem recall that \mathcal{P} is the set of transition probabilities parameterized by $\theta \in \Theta$. We index time $t = 1, 2, \dots$ in a continuous fashion. Episode $k = 1, 2, \dots$ and stage $h = 1, \dots, H - 1$ corresponds to time $t = (k - 1)(H - 1) + h$:

episode (k)	1	1	...	1	2	2	...	2	3	...
stage (h)	1	2	...	$H - 1$	1	2	...	$H - 1$	1	...
time step (t)	1	2	...	$H - 1$	H	$H + 1$...	$2H - 2$	$2H - 1$...

Note that the transitions at stage $h = H$ are skipped and the time index at the end of episode $k \geq 1$ is $k(H - 1)$.

Let $V_{(t)}$ be the value function used by Algorithm 1 at time t ($V_{(t)}$ is constant in periods of length $H - 1$), while let $(s_{(t)}, a_{(t)})$ be the state-action pair visited at time t .

Let \mathcal{V} be the set of optimal value functions under some model in \mathcal{P} : $\mathcal{V} = \{V_{P'}^* : P' \in \mathcal{P}\}$. Note that $\mathcal{V} \subset \mathcal{B}(\mathcal{S}, H)$, where $\mathcal{B}(\mathcal{S}, H)$ denotes the set of real-valued measurable functions with domain \mathcal{S} that are bounded by H . Note also that for all t , $V_{(t)} \in \mathcal{V}$. Define $\mathcal{X} = \mathcal{S} \times \mathcal{A} \times \mathcal{V}$. We also let $X_t = (s_{(t)}, a_{(t)}, V_{(t)})$, $Y_t = V_{(t)}(s_{(t+1)})$ when $t + 1 \notin \{H + 1, 2H + 1, \dots\}$ and $Y_t = V_{(t)}(s_{H+1}^k)$, and choose

$$\mathcal{F} = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : \exists \tilde{P} \in \mathcal{P} \text{ s.t. } f(s, a, v) = \int \tilde{P}_a(ds'|s)v(s') \right\}. \quad (11)$$

Note that $\mathcal{F} \subset \mathcal{B}_\infty(\mathcal{X}, H)$.

Let $\phi : \mathcal{P} \rightarrow \mathcal{F}$ be the natural surjection to \mathcal{F} : $\phi(P) = f$ where $f(s, a, v) = \int P_a(ds'|s)v(s')$ for $(s, a, v) \in \mathcal{X}$. We know show that ϕ is in fact a bijection. If $P \neq P'$, this means that for some $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $U \subset \mathcal{S}$ measurable, $P_a(U|s) \neq P'_a(U|s)$. Choosing v to be the indicator of U , note that $(s, a, v) \in \mathcal{X}$. Hence, $\phi(P)(s, a, v) = P_a(U|s) \neq$

$P'_a(U|s) = \phi(P')(s, a, v)$, and hence $\phi(P) \neq \phi(P')$: ϕ is indeed a bijection. For convenience and to reduce clutter, we will write $f_P = \phi(P)$.

Choose $\mathbb{F} = (\mathbb{F}_t)_{t \geq 0}$ so that \mathbb{F}_{t-1} is generated by $(s_{(1)}, a_{(1)}, V_{(1)}, \dots, s_{(t)}, a_{(t)}, V_{(t)})$. Then $\mathbb{E}[Y_t | \mathbb{F}_{t-1}] = \int P_{a_{(t)}}(ds' | s_{(t)}) V_{(t)}(s') = f_P(X_t)$ and by definition $f_P \in \mathcal{F}$. Now, $Y_t \in [0, H]$, hence, $Z_t = Y_t - f_P(X_t)$ is conditionally $H/2$ -subgaussian given \mathbb{F}_{t-1} .

Let $t = k(H-1)$ for some $k \geq 1$. Thus, this time step corresponds to finishing episode k and thus $V_{(t)} = V_k$. Furthermore, letting $\hat{f}_t = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{p=1}^t (f(X_p) - Y_p)^2$, since ϕ is an injection, we see that $\hat{f}_t = f_{\hat{P}_k}$ where \hat{P}_k is defined using (9). For $P', P'' \in \mathcal{P}$, we have $L_k(P', P'') = \sum_{p=1}^t (f_{P'}(X_p) - f_{P''}(X_p))^2$ and thus

$$\begin{aligned} \mathcal{B}_k &= \{\tilde{P} \in \mathcal{P} : L_k(\hat{P}_k, \tilde{P}) \leq \beta_k\} = \{\tilde{P} \in \mathcal{P} : \sum_{p=1}^t (\hat{f}_t(X_p) - f_{\tilde{P}}(X_p))^2 \leq \beta_k\} \\ &= \{\phi^{-1}(f) : f \in \mathcal{F} \text{ and } \sum_{p=1}^t (\hat{f}_t(X_p) - f(X_p))^2 \leq \beta_k\} = \phi^{-1}(\mathcal{F}_t(\beta_k)). \end{aligned}$$

Corollary 7. For $\alpha > 0$ and $k \geq 1$ let

$$\beta_k = 2H^2 \log \left(\frac{2\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) + 2H(kH-1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{4kH(kH-1)}{\delta} \right)} \right\}.$$

Then, with probability $1 - \delta$, for any $k \geq 1$, $P \in \mathcal{B}_k$ where \mathcal{B}_k is defined by (10).

A.5. Regret of Algorithm 1

Recall that $\mathcal{X} = \mathcal{S} \times \mathcal{A} \times \mathcal{V}$ where $\mathcal{V} \subset \mathcal{B}_\infty(\mathcal{S}, H)$ is the set of value functions that are optimal under some model in \mathcal{P} . We will abbreviate $(x_1, \dots, x_t) \in \mathcal{X}^t$ as $x_{1:t}$. Further, we let $\mathcal{F}|_{x_{1:t}} = \{(f(x_1), \dots, f(x_t)) : f \in \mathcal{F}\} \subset \mathbb{R}^t$ and for $S \subset \mathbb{R}^t$, let $\operatorname{diam}(S) = \sup_{u, v \in S} \|u - v\|_2$ be the diameter of S . We will need the following lemma, extracted from [Russo & Van Roy \(2014\)](#):

Lemma 8 (Lemma 5 of [Russo & Van Roy \(2014\)](#)). *Let $\mathcal{F} \subset \mathcal{B}_\infty(\mathcal{X}, C)$ be a set of functions bounded by $C > 0$, $(\mathcal{F}_t)_{t \geq 1}$ and $(x_t)_{t \geq 1}$ be sequences such that $\mathcal{F}_t \subset \mathcal{F}$ and $x_t \in \mathcal{X}$ hold for $t \geq 1$. Then, for any $T \geq 1$ and $\alpha > 0$ it holds that*

$$\sum_{t=1}^T \operatorname{diam}(\mathcal{F}_t|_{x_{1:t}}) \leq \alpha + C(d \wedge T) + 2\delta_T \sqrt{dT},$$

where $\delta_T = \max_{1 \leq t \leq T} \operatorname{diam}(\mathcal{F}_t|_{x_{1:t}})$ and $d = \dim_{\mathcal{E}}(\mathcal{F}, \alpha)$.

Let

$$W_k = \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}_{a_h^k}(s_h^k) - P_{a_h^k}(s_h^k), V_{h,k} \rangle.$$

From Lemma 5, we get

$$R_K \leq \sum_{k=1}^K W_k + \sum_{k=1}^K \sum_{h=1}^{H-1} \xi_{h+1,k}. \quad (12)$$

Lemma 9. *Let $\alpha > 0$ and $d = \dim_{\mathcal{E}}(\mathcal{F}, \alpha)$ where \mathcal{F} is given by (11). Then, for any nondecreasing sequence $(\beta_k^2)_{k=1}^K$, on the event when $P \in \cap_{k \in [K]} \mathcal{B}_k$,*

$$\sum_{k=1}^K W_k \leq \alpha + H(d \wedge K(H-1)) + 4\sqrt{d\beta_K K(H-1)}.$$

Proof. Let $P \in \bigcap_{k \in [K]} \mathcal{B}_k$ holds. Using the notation of the previous section, letting $\tilde{\mathcal{F}}_t = \mathcal{F}_t(\beta_k)$ for $(k-1)(H-1) + 1 \leq t \leq k(H-1)$, we have

$$\begin{aligned} \sum_{k=1}^K W_k &\leq \sum_{k=1}^K \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} (f_{\tilde{P}}(s_h^k, a_h^k, V_{h+1,k}) - f_P(s_h^k, a_h^k, V_{h+1,k})) \\ &\leq \sum_{t=1}^{K(H-1)} \text{diam}(\tilde{\mathcal{F}}_t|_{X_t}) \quad (\text{because } P \in \bigcap_{k \in [K]} \mathcal{B}_k) \\ &\leq \alpha + H(d \wedge K(H-1)) + 2\delta_{K(H-1)} \sqrt{dK(H-1)}, \end{aligned}$$

where X_t is defined in Section A.4 and where the last inequality is by Lemma 8, which is applicable because $\mathcal{F} \subset \mathcal{B}_\infty(\mathcal{X}, H)$ holds by choice, and $\delta_{K(H-1)} = \max_{1 \leq t \leq K(H-1)} \text{diam}(\tilde{\mathcal{F}}_t|_{X_{1:t}})$. Thanks to the definition of $\tilde{\mathcal{F}}_t$, $\delta_{K(H-1)} \leq 2\sqrt{\beta_K}$. Plugging this into the previous display finishes the proof. \square

A.5.1. PROOF OF THEOREM 1

Proof. Note that for any $k \in [K]$ and $h \in [H-1]$, $\xi_{h+1,k} \in [-H, H]$. As noted beforehand, $\xi_{2,1}, \xi_{3,1}, \dots, \xi_{H,1}, \xi_{2,2}, \xi_{3,2}, \dots, \xi_{H,2}, \xi_{2,3}, \dots$ is a martingale difference sequence. Thus, with probability $1 - \delta$, $\sum_{k=1}^K \sum_{h=1}^{H-1} \xi_{h+1,k} \leq H\sqrt{2K(H-1)\log(1/\delta)}$. Consider the event when this inequality holds and when $P \in \bigcap_{k \in [K]} \mathcal{B}_k$. By using Corollary 7 and a union bound, this event holds with probability at least $1 - 2\delta$. On this event, by (12) and Lemma 9, we obtain

$$R_K \leq \alpha + H(d \wedge K(H-1)) + 4\sqrt{d\beta_K K(H-1)} + H\sqrt{2K(H-1)\log(1/\delta)}.$$

Using $\alpha \leq 1$, which holds by assumption, finishes the proof. \square

A.5.2. PROOF OF COROLLARY 2

Proof. Note that

$$\begin{aligned} \|f_{P'} - f_{P''}\|_\infty &= \sup_{s,a,v} \left| \int (P'_a(ds'|s) - P''_a(ds'|s))v(s') \right| \leq H \sup_{s,a} \int |P'_a(ds'|s) - P''_a(ds'|s)| \\ &= H \sup_{s,a} \|P'_a(s) - P''_a(s)\|_1 =: H\|P' - P''\|_{\infty,1}. \end{aligned}$$

For $\alpha > 0$ let $\mathcal{N}(\mathcal{P}, \alpha, \|\cdot\|_{\infty,1})$ denote the $(\alpha, \|\cdot\|_{\infty,1})$ -covering number of \mathcal{P} . Then we have

$$\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty) \leq \mathcal{N}(\mathcal{P}, \alpha/H, \|\cdot\|_{\infty,1}).$$

Then, by Corollary 7,

$$\beta_K = 2H^2 \log(2\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)/\delta) + C \leq 2H^2 \log(2\mathcal{N}(\mathcal{P}, \alpha/H, \|\cdot\|_{\infty,1})/\delta) + C$$

with some universal constant $C > 0$. Let $f : (\Theta, \|\cdot\|) \rightarrow (\mathcal{P}, \|\cdot\|_{\infty,1})$ be defined by $\theta \mapsto \sum_j \theta_j P_j$. Note that $\|f(\theta) - f(\theta')\|_{\infty,1} \leq \sup_{s,a} \sum_j \|(\theta_j - \theta'_j)P_{j,a}(s)\|_1 = \sum_j |\theta_j - \theta'_j| = \|\theta - \theta'\|_1$. Hence, any $(\epsilon, \|\cdot\|_1)$ covering of Θ induces an $(\epsilon, \|\cdot\|_{\infty,1})$ -covering of \mathcal{P} and so $\mathcal{N}(\mathcal{P}, \alpha/H, \|\cdot\|_{\infty,1}) \leq \mathcal{N}(\Theta, \alpha/H, \|\cdot\|_1) \leq C'(RH/\alpha)^d$ with some universal constant $C' > 0$.

Now, choose $1/\alpha = K\sqrt{\log(KH/\delta)}$. Hence,

$$\beta_K \leq 2H^2(\log(2C'/\delta) + d \log(RH/\alpha)) + C.$$

Suppressing log factors (e.g., $\log(RH)$), log log terms and constants, we have $\beta_K = H^2(d + \log(1/\delta))$.

Let \mathcal{F} be given by (11). We now bound $\dim_\mathcal{E}(\mathcal{F}, \alpha)$. Let $\mathcal{X} = \mathcal{S} \times \mathcal{A} \times B(\mathcal{S})$ as before. Define $z : \mathcal{S} \times \mathcal{A} \times B(\mathcal{S}) \rightarrow \mathbb{R}^d$ using $z(s, a, v)_j = \langle P_{j,a}(s), v \rangle$ and note that if $x \in \mathcal{X}$ is (ϵ, \mathcal{F}) -independent of $x_1, \dots, x_k \in \mathcal{X}$ then $z(x) \in \mathbb{R}^d$ is (ϵ, Θ) -independent of $z(x_1), \dots, z(x_k) \in \mathbb{R}^d$. This holds because if $P = \sum_j \theta_j P_j \in \mathcal{P}$ then $f_P(s, a, v) = \langle \theta, z(s, a, v) \rangle$

for any $(s, a, v) \in \mathcal{X}$. Hence, $\dim_{\mathcal{E}}(\mathcal{F}, \alpha) \leq \dim_{\mathcal{E}}(\text{Lin}(\mathcal{Z}, \Theta), \alpha)$, where $\text{Lin}(\mathcal{Z}, \Theta)$ is the set of linear maps with domain $\mathcal{Z} = \{z(x) : x \in \mathcal{X}\} \subset \mathbb{R}^d$ and parameter from Θ : $\text{Lin}(\mathcal{Z}, \Theta) = \{h : h : \mathcal{Z} \rightarrow \mathbb{R} \text{ s.t. } \exists \theta \in \Theta : h(z) = \langle \theta, z \rangle, z \in \mathcal{Z}\}$. Now, by Proposition 11 of Russo & Van Roy (2014), $\dim_{\mathcal{E}}(\text{Lin}(\mathcal{Z}, \Theta), \alpha) = O(d \log(1 + (S\gamma/\alpha)^2))$ where S is the $\|\cdot\|_2$ diameter of Θ and $\gamma = \sup_{z \in \mathcal{Z}} \|z\|_2$. We have

$$\|z\|_2^2 = \sum_j (\langle P_{j,a}(s), v \rangle)^2 \leq H^2 d,$$

hence $\gamma \leq H\sqrt{d}$. By the relation between the 1 and 2 norms, the 2-norm diameter of Θ is at most $\sqrt{d}R$. Dropping log terms, $\dim_{\mathcal{E}}(\mathcal{F}, \alpha) = \tilde{O}(d)$.

Plugging into Theorem 1 gives the desired result. \square

A.6. Proof of Theorem 6

Recall the following:

Definition 3. A random variable X is σ -subgaussian if for all $\lambda \in \mathbb{R}$, it holds that $\mathbb{E}[\exp(\lambda X)] \leq \exp(\lambda^2 \sigma^2 / 2)$.

The proof of the next couple of statements is standard and is included only for completeness.

Theorem 10. If X is σ -subgaussian, then for any $\lambda > 0$, with probability at least $1 - \delta$,

$$X < \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right) + \lambda \frac{\sigma^2}{2}. \quad (13)$$

Proof. Let $\lambda > 0$. We have, $\{X \geq \epsilon\} = \{\exp(\lambda(X - \epsilon)) \geq 0\}$. Hence, Markov's inequality gives $\mathbb{P}(X \geq \epsilon) \leq \exp(-\lambda\epsilon) \mathbb{E}[\exp(\lambda X)] \leq \exp(-\lambda\epsilon + \frac{1}{2}\lambda^2\sigma^2)$. Equating the right-hand side with δ and solving for ϵ , we get that $\log(\delta) = -\lambda\epsilon + \frac{1}{2}\lambda^2\sigma^2$. Solving for ϵ gives $\epsilon = \log(1/\delta)/\lambda + \frac{\sigma^2}{2}\lambda$, finishing the proof. \square

Choosing the λ that minimizes the right-hand side of the bound gives the usual form:

$$\mathbb{P}(X \geq \sqrt{2\sigma^2 \log(1/\delta)}) \leq \delta. \quad (14)$$

Lemma 11 (Lemma 5.4 of Lattimore & Szepesvári (2020)). Suppose that X is σ -subgaussian and X_1 and X_2 are independent and σ_1 and σ_2 -subgaussian, respectively, then:

1. $\mathbb{E}[X] = 0$.
2. cX is $|c|\sigma$ -subgaussian for all $c \in \mathbb{R}$.
3. $X_1 + X_2$ is $\sqrt{\sigma_1^2 + \sigma_2^2}$ -subgaussian.

Let $(Z_p)_p$ be an $\mathbb{F} = (\mathbb{F}_p)_p$ -adapted process. Recall that $(Z_p)_p$ is conditionally σ -subgaussian given \mathbb{F} if for all $p \geq 1$,

$$\log \mathbb{E}[\exp(\lambda Z_p) | \mathbb{F}_{p-1}] \leq \frac{1}{2} \lambda^2 \sigma^2, \quad \text{for all } \lambda \in \mathbb{R}.$$

A standard calculation gives that $S_t = \sum_{p=1}^t Z_p$ is $\sqrt{t}\sigma$ -subgaussian (essentially, a refinement of the calculation that is need to show Part (3) of Lemma 11) and thus, in particular, for any $t \geq 1$ and $\lambda > 0$, with probability $1 - \delta$,

$$S_t < \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right) + \lambda \frac{t\sigma^2}{2}.$$

In fact, by slightly strengthening the argument, one can show that the above inequality holds simultaneously for all $t \geq 1$:

Theorem 12 (E.g., Lemma 7 of Russo & Van Roy (2014)). Let \mathbb{F} be a filtration and let $(Z_p)_p$ be an \mathbb{F} -adapted, conditionally σ -subgaussian process. Then for any $\lambda > 0$, with probability at least $1 - \delta$, for all $t \geq 1$,

$$S_t < \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right) + \lambda \frac{t\sigma^2}{2}, \quad (15)$$

where $S_t = \sum_{p=1}^t Z_p$.

Proof of Theorem 6 Let us introduce the following helpful notation: For vectors $x, y \in \mathbb{R}^t$, let $\langle x, y \rangle_t = \sum_{p=1}^t x_p y_p$, $\|x\|_t^2 = \langle x, x \rangle_t$, and for $f : \mathcal{X} \rightarrow \mathbb{R}$, $\|f\|_t^2 = \sum_{p=1}^t f^2(X_p)$. More generally, we will overload addition and subtraction such that for $x \in \mathbb{R}^t$, $x + f \in \mathbb{R}^t$ is the vector whose p th coordinate is $x_p + f(X_p)$ (x_p and X_p both appear on purpose here). We also overload $\langle \cdot, \cdot \rangle_t$ such that $\langle x, f \rangle_t = \langle f, x \rangle_t = \sum_{p=1}^t x_p f(X_p)$.

Define Z_p using $Y_p = f_*(X_p) + Z_p$ and collect $(Y_p)_{p=1}^t$ and $(Z_p)_{p=1}^t$ into the vectors Y and Z . As in the statement of the theorem, let $\mathbb{F} = (\mathbb{F}_p)_{p=0,1,\dots}$ be such that for any $s \geq 1$, $(X_1, Y_1, \dots, X_{p-1}, Y_{p-1}, X_p)$ is \mathbb{F}_{p-1} -measurable. Note that for any $p \geq 1$, $Z_p = Y_p - f_*(X_p)$ is \mathbb{F}_p -measurable, hence $(Z_p)_{p \geq 1}$ is \mathbb{F} -adapted.

With this, elementary calculation gives

$$\|Y - f\|_t^2 - \|Y - f_*\|_t^2 = \|f_* - f\|_t^2 + 2\langle Z, f_* - f \rangle_t.$$

Splitting $\|f_* - f\|_t^2$ and rearranging gives

$$\frac{1}{2}\|f_* - f\|_t^2 = \|Y - f\|_t^2 - \|Y - f_*\|_t^2 + E(f) \quad (16)$$

where

$$E(f) = -\frac{1}{2}\|f_* - f\|_t^2 + 2\langle Z, f - f_* \rangle_t.$$

Recall that $\hat{f}_t = \operatorname{argmin}_{f \in \mathcal{F}} \|Y - f\|_t^2$. Plugging \hat{f}_t into 16 in place of f and using that thanks to $f_* \in \mathcal{F}$, $\|Y - \hat{f}_t\|_t^2 \leq \|Y - f_*\|_t^2$, we get

$$\frac{1}{2}\|f_* - \hat{f}_t\|_t^2 \leq E(\hat{f}_t). \quad (17)$$

Thus, it remains to bound $E(\hat{f}_t)$. For this fix some $\alpha > 0$ to be chosen later and let $\mathcal{G}(\alpha) \subset \mathcal{F}$ be an α -cover of \mathcal{F} in $\|\cdot\|_\infty$. Let $g \in \mathcal{G}(\alpha)$ be a random function, also to be chosen later. We have

$$E(\hat{f}_t) = E(\hat{f}_t) - E(g) + E(g) \leq E(\hat{f}_t) - E(g) + \max_{\tilde{g} \in \mathcal{G}(\alpha)} E(\tilde{g}) \quad (18)$$

We start by bounding the last term above. A simple calculation gives that for any fixed $f \in \mathcal{F}$, w.p. $1 - \delta$, $2\langle Z, f - f_* \rangle_t$ is $2\sigma\|f - f_*\|_t$ -subgaussian. Hence, with probability $1 - \delta$, simultaneously for all $t \geq 1$,

$$E(f) \leq -\frac{1}{2}\|f_* - f\|_t^2 + \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right) + \lambda \frac{4\sigma^2\|f - f_*\|_t^2}{2} = 4\sigma^2 \log\left(\frac{1}{\delta}\right),$$

where the equality follows by choosing $\lambda = 1/(4\sigma^2)$ (which makes the first and last terms cancel). (Note how splitting $\|f - f_*\|_t^2$ into two halves allowed us to bound the ‘‘error term’’ $E(f)$ independently of t .) Now, by a union bound, it follows that with probability at least $1 - \delta$, the second term is bounded by $4\sigma^2 \log(|\mathcal{G}(\alpha)|/\delta)$.

Let us now turn to bounding the first term. We calculate

$$\begin{aligned} E(\hat{f}_t) - E(g) &= \frac{1}{2}\|g - f_*\|_t^2 - \frac{1}{2}\|\hat{f}_t - f_*\|_t^2 + 2\langle Z, \hat{f}_t - g \rangle_t \\ &\leq \frac{1}{2} \left(\langle g - \hat{f}_t, g + \hat{f}_t + 2f_* \rangle_t \right) + 2\|Z\|_t \|\hat{f}_t - g\|_t \\ &\leq \frac{1}{2} 4C\alpha t + 2\|Z\|_t \alpha \sqrt{t}, \end{aligned}$$

where for the last inequality we chose $g = \operatorname{argmin}_{\tilde{g} \in \mathcal{G}(\alpha)} \|\hat{f}_t - \tilde{g}\|_\infty$ so that $\|\hat{f}_t - g\|_t \leq \alpha\sqrt{t}$ and used Cauchy-Schwartz, together with that $\|g\|_t, \|\hat{f}_t\|_t, \|f_*\|_t \leq C\sqrt{t}$, which follows from $g, \hat{f}_t, f_* \in \mathcal{F}$ and that by assumption all functions in \mathcal{F} are bounded by C .

It remains to bound $\|Z\|_t$. For this, we observe that with probability $1 - \delta$, simultaneously for all $t \geq 1$,

$$\|Z\|_t \leq \sigma \sqrt{2t \log(2t(t+1)/\delta)}.$$

Indeed, this follows because with probability $1 - \delta$, simultaneously for any $s \geq 1$, $|Z_p|^2 \leq 2\sigma^2 \log(2s(s+1)/\delta)$ holds because of a union bound and Eq. (14). Therefore, for the above choice g , with probability $1 - \delta$, simultaneously for all $t \geq 1$, it holds that

$$E(\hat{f}_t) - E(g) \leq 2C\alpha t + 2t\alpha \sqrt{\sigma^2 \log(2t(t+1)/\delta)}.$$

Merging this with Eqs. (17) and (18) and with another union bound, we get that with probability $1 - \delta$, for any $t \geq 1$,

$$\|f_* - \hat{f}_t\|_t^2 \leq 8\sigma^2 \log(2N_\alpha/\delta) + 4t\alpha \left(C + \sqrt{\sigma^2 \log(4t(t+1)/\delta)} \right),$$

where N_α is the $(\alpha, \|\cdot\|_\infty)$ -covering number of \mathcal{F} . □

B. Proof of Theorem 3

In this section we establish a regret lower bound by reduction to a known result for tabular MDP.

Proof. We assume without loss of generality that d is a multiple of 4 and $d \geq 8$. We set $S = 2$ and $A = d/4 \geq 2$. According to (Azar et al., 2017), (Osband & Van Roy, 2016), there exists an MDP $\mathcal{M}(\mathcal{S}, \mathcal{A}, P, r, H)$ with S states, A actions and horizon H such that any algorithm has regret at least $\Omega(\sqrt{HSAT})$. In this case, we have $|\mathcal{S} \times \mathcal{A} \times \mathcal{S}| = d$. We use $\sigma(s, a, s')$ to denote the index of (s, a, s') in $\mathcal{S} \times \mathcal{A} \times \mathcal{S}$. Letting

$$P_i(s'|s, a) = \begin{cases} 1 & \text{if } \sigma(s, a, s') = i, \\ 0 & \text{otherwise,} \end{cases}$$

and $\theta^i = P(s'|s, a)$ if $\sigma(s, a, s') = i$, we will have $P(s'|s, a) = \sum_{i=1}^d \theta^i P_i(s'|s, a)$. Therefore P can be parametrized using (19). Therefore, the known lower bound $\Omega(\sqrt{HSAT})$ implies a worst-case lower bound of $\Omega(\sqrt{H \cdot d/2 \cdot T}) = \Omega(\sqrt{HdT})$ for our model. □

C. The Special Case of Linear Transition Models

We derive a modification of UCRL-VTR when P_θ is a linear model of the form $P_\theta = \sum_{j=1}^d \theta_j P_j$, which is captured in the following assumption:

Assumption 3 (Linear Parameterized Transition Model). *There exists a vector $\theta_* \in \mathbb{R}^d$ such that $\|\theta_*\|_2 \leq C_\theta$ ($C_\theta \geq 1$) and*

$$P(s'|s, a) = \sum_{j=1}^d (\theta_*)_j P_j(s'|s, a) = P.(s'|s, a)^\top \theta_*, \quad (19)$$

where P_j 's are known basis models such that $\sup_{j \in [d], (s,a) \in \mathcal{S} \times \mathcal{A}} \|P_j(\cdot|s, a)\|_1 \leq 1$, and $P.(s'|s, a)$ denotes the d -dimensional vector $P.(s'|s, a) = [P_1(s'|s, a), \dots, P_d(s'|s, a)]^\top$ ³. Note that we do not require each basis model P_j to be a probability transition model.

By modifying the algorithm and using optimistic Q-update, we obtain an algorithm that can be implemented using efficient recursive update. See Algorithm 3 for full details of implementation.

Estimating θ_* by recursive regression. We let $X_{h,k}^\top \theta := \mathbb{E}.[V_{h+1,k}(s)|s_h^k, a_h^k]^\top \theta = \langle P_\theta(\cdot|s, a), V_{h+1,k} \rangle$ be the predicted expected value of next state. In this case, each new observation adds the following loss to regression:

$$\begin{aligned} & (X_{h,k}^\top \theta - y_{h,k})^2 : \\ & = (\mathbb{E}.[V_{h+1,k}(s)|s_h^k, a_h^k]^\top \theta - V_{h+1,k}(s_{h+1}^k))^2 \end{aligned}$$

³We also use $P.(\cdot|s, a)$ to denote a $d \times S$ matrix.

Algorithm 3 UCRL-VTR with linear transition model

- 1: **Input:** MDP, $d, H, T = KH$;
- 2: **Initialize:** $M_{1,1} \leftarrow H^2 dI$, $w_{1,1} \leftarrow 0 \in \mathbb{R}^{d \times 1}$, $\theta_1 \leftarrow M_{1,1}^{-1} w_{1,1}$ for $1 \leq h \leq H$;
- 3: **Initialize:** $\delta \leftarrow 1/K$, and for $1 \leq k \leq K$,

$$\beta_k \leftarrow 16C_\theta^2 H^2 d \log(1 + Hk) \log^2((k+1)^2 H / \delta);$$

- 4: Compute Q-function $Q_{h,1}$ using $\theta_{1,1}$ according to (3);

5: **for** $k = 1 : K$ **do**

- 6: Obtain initial state s_1^k for episode k ;

7: **for** $h = 1 : H$ **do**

- 8: Choose action greedily by

$$a_h^k = \arg \max_{a \in \mathcal{A}} Q_{h,k}(s_h^k, a)$$

and observe the next state s_{h+1}^k .

- 9: Compute the predicted value vector:

▷ Evaluate the expected value of next state

$$\begin{aligned} X_{h,k} &\leftarrow \mathbb{E}.[V_{h+1,k}(s) | s_h^k, a_h^k] \\ &= \sum_{s \in \mathcal{S}} V_{h+1,k}(s) \cdot P.(s | s_h^k, a_h^k). \end{aligned}$$

- 10: $y_{h,k} \leftarrow V_{h+1,k}(s_{h+1}^k)$

▷ Update regression parameters

- 11: $M_{h+1,k} \leftarrow M_{h,k} + X_{h,k} X_{h,k}^\top$

- 12: $w_{h+1,k} \leftarrow w_{h,k} + y_{h,k} \cdot X_{h,k}$

13: **end for**

- 14: Update at the end of episode:

▷ Update Model Parameters

$$M_{1,k+1} \leftarrow M_{H+1,k},$$

$$w_{1,k+1} \leftarrow w_{H+1,k},$$

$$\theta_{k+1} \leftarrow M_{1,k+1}^{-1} w_{1,k+1};$$

- 15: Compute $Q_{h,k+1}$, $h = H, \dots, 1$, using θ_{k+1} according to (20)

▷ Computing Q functions

16: **end for**

By aggregating the value prediction losses constructed from all past experiences, we formulate a ridge regression problem to estimate θ_* by

$$\begin{aligned} & \theta_{k+1} \\ &= \arg \min_{\theta \in \mathbb{R}^d} \left[\theta^\top M_{1,1} \theta + \sum_{(h',k') \leq (H,k)} (X_{h',k'}^\top \theta - y_{k',h'})^2 \right], \end{aligned}$$

where $M_{1,1} = H^2 d I$ acts as a regularization term.

To solve the above regression problem, we can first calculate $X_{h',k'}$ and recursively compute estimates of θ_* by letting

$$\begin{aligned} M_{1,k+1} &= M_{1,1} + \sum_{(h',k') \leq (H,k)} X_{h',k'} X_{h',k'}^\top \\ w_{1,k+1} &= w_{1,1} + \sum_{(h',k') \leq (H,k)} y_{h',k'} \cdot X_{h',k'}, \end{aligned}$$

with $M_{1,1} = H^2 d \cdot I$ and $w_{1,1} = 0$. Then we obtain the estimated θ_{k+1} easily by

$$\theta_{k+1} = M_{1,k+1}^{-1} w_{k+1}.$$

Confidence ball. We construct B_k as follows:

$$B_k = \{\theta | (\theta - \theta_k)^\top M_k (\theta - \theta_k) \leq \beta_k\}.$$

where β_k is preselected (see the algorithm).

Our model parameter update, θ_k and M_k , can be via a recursive update in an incremental fashion. In this way, one does not need to re-train the model parameter from scratch every episode. A similarly simple recursion was used in (Jin et al., 2019) for model-free Q learning. Our method differs in that our Q functions cannot be parameterized by d parameters and our updates are made on the transition model rather than Q functions.

Optimistic Q-update. Instead of solving the optimistic planning problem $\theta_k = \arg \max_{\theta} \{V_{\theta}^*(s_1) | \theta \in B_k\}$ as in Algorithm 1, we incorporate optimism into iterative Q-update:

$$\begin{aligned} Q_{H+1,k}(s, a) &= 0, \\ V_{h,k}(s) &= \max_{a \in \mathcal{A}} Q_{h,k}(s, a), \\ Q_{h,k}(s, a) &= r(s, a) + \max_{\theta \in B_k} \sum_{j=1}^d (\theta)_j P_j(\cdot | s, a) V_{h+1,k}. \end{aligned}$$

Since the confidence sets are ellipsoids, the preceding Q update has a closed-forms solution

$$\begin{aligned} & Q_{h,k}(s, a) \\ &= r(s, a) + \max_{\theta \in B_k} \langle P_{\theta}(\cdot | s, a), V_{h+1,k} \rangle \\ &= r(s, a) + X_{h,k}^\top \theta_k + \sqrt{\beta_k} \sqrt{X_{h,k}^\top M_k^{-1} X_{h,k}}. \end{aligned} \tag{20}$$

The last term in the above is the ‘‘bonus’’ term that quantifies uncertainty and encourages exploration. This optimistic Q value allows us to greedily pick actions while sufficiently exploring the state space.

Algorithm 3 is a modification of UCRL-VTR and uses a different construction of confidence set. we provide an independent regret analysis using techniques from linear bandit theory. The next theorem gives a regret upper bound for Algorithm 3.

Theorem 13. *Let Assumption 3 hold. If we choose*

$$\beta_k = \left(H \sqrt{d \log \left(\frac{1 + Hk \cdot H^2 d}{\delta} \right)} + C_{\theta} H \sqrt{d} \right)^2,$$

then T -time-step regret of Algorithm 1 satisfies

$$\mathbb{E}[R(T)] = \tilde{O}\left(C_\theta \cdot d\sqrt{H^3 T}\right),$$

where $T = HK$ is the total number of steps in K episodes, C_θ ($C_\theta \geq 1$) is a known constant such that $\|\theta_*\| \leq C_\theta$ and \tilde{O} hides polylog factors of H, T .

Let us outline the proof ideas. In the first part of the proof, we show that if $\theta_* \in B_{h,k}$, then the estimated Q-functions are optimistic estimates of the true Q-value functions. That is, $Q_{h,k}(s)$ is greater than the true Q-value $Q_h(s)$ for every $s \in \mathcal{S}$. Using this fact, we can bound the regret by the sum of $Q_{1,k}(s_1^k) - Q_1^{\pi_k}(s_1^k)$, which can be decomposed into the sum of state-action confidence bounds on the sample path. In the second part, we construct martingale difference sequences and apply a concentration argument to show that $\theta_* \in B_{h,k}$ for all (h, k) with high probability. The full proof is deferred to the Appendix E.

D. Proof of Theorem 4

In this section, we will present the full proof of Theorem 4. To handle the misspecification error, we will modify the bonus term by replacing it with

$$\beta_k = 8H^2 \log\left(\frac{2\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta}\right) + 4H(kH-1)\alpha \left\{ 2 + \sqrt{\log\left(\frac{4kH(kH-1)}{\delta}\right)} \right\} + 8H^3 k \varepsilon^2.$$

The last term in the above choice of β_k can be viewed as an ‘‘error tolerance.’’

Next we show that $P^* \in B_k$ with high probability.

We first present a theorem which is nearly identical to Theorem 6 but tolerates misspecification. We use the same notations as in the proof of Theorem 6.

Theorem 14. *Let \mathbb{F} be the filtration defined above and assume that the functions in \mathcal{F} and also f_* are all bounded by the positive constant $C > 0$ at values X_t for all t . Assume that there exists $\tilde{f} \in \mathcal{F}$ such that $|\tilde{f}(X) - f_*(X)| \leq \zeta$ for all $X = (s, a, v)$ with $\|v\|_\infty \leq H$, and also for each $s \geq 1$, $(Y_p - f_*(X_p))_p$ is conditionally σ -subgaussian given \mathbb{F}_{p-1} . We define*

$$\hat{f}_t = \arg \min_{f \in \mathcal{F}} \sum_{p=1}^t (f(X_p) - Y_p)^2$$

and

$$\mathcal{F}_t(\beta) = \left\{ f : \mathcal{X} \rightarrow \mathbb{R}, \text{ s.t. } \sum_{p=1}^t (f(X_p) - \hat{f}_t(X_p))^2 \leq \beta \right\}.$$

Then, for any $\alpha > 0$, with probability $1 - \delta$, for all $t \geq 1$, $f_* \in \mathcal{F}_t(\beta_t(\delta, \alpha))$, where

$$\beta_t(\delta, \alpha) = 16\sigma^2 \log(4N_\alpha/\delta) + 4t\alpha \left(C + \sqrt{\sigma^2 \log(8t(t+1)/\delta)} \right) + 3t\zeta^2.$$

Note that here the last term is due to the misspecification error.

Proof. The proof of this theorem is also nearly identical to Theorem 6, except for the modifications below.

Due to model misspecification, we no longer have $f_* \in \mathcal{F}$, and hence we may not have $\|Y - \hat{f}_t\|_t \leq \|Y - f_*\|_t$ (Here notation $\|\cdot\|_t$ is defined to be the same as the notations in Theorem 6). To handle the misspecification error, we will use the function \tilde{f} as a bridge to bound the error between \hat{f}_t and f_* . Hence since $\hat{f}_t = \arg \min_{f \in \mathcal{F}} \|Y - f\|_t^2$, we have $\|\hat{f}_t - Y\|_t^2 \leq \|\tilde{f} - Y\|_t^2$, which indicates that $\|\hat{f}_t - f_* - Z\|_t^2 \leq \|\tilde{f} - f_* - Z\|_t^2$. (Recall the notations $Z_p = Y_p - f_*(X_p)$ and $Z = (Z_1, \dots, Z_p)$.) Therefore, we have

$$\|\hat{f}_t - f_*\|_t^2 - 2\langle \hat{f}_t - f_*, Z \rangle_t \leq \|\tilde{f} - f_*\|_t^2 - 2\langle \tilde{f} - f_*, Z \rangle_t.$$

We then obtain

$$\begin{aligned} \frac{1}{2}\|\hat{f}_t - f_*\|_t^2 &\leq -\frac{1}{2}\|\hat{f}_t - f_*\|_t^2 + 2\langle \hat{f}_t - f_*, Z \rangle_t + \|\tilde{f} - f_*\|_t^2 - 2\langle \tilde{f} - f_*, Z \rangle_t \\ &= E(\hat{f}_t) + \tilde{E}(\tilde{f}) + \frac{3}{2}\|\tilde{f} - f_*\|_t^2, \end{aligned} \quad (21)$$

where we define

$$E(f) = -\frac{1}{2}\|f - f_*\|_t^2 + 2\langle Z, f - f_* \rangle_t \quad (22)$$

$$\tilde{E}(f) = -\frac{1}{2}\|f - f_*\|_t^2 - 2\langle Z, f - f_* \rangle_t \quad (23)$$

In the next, we will bound $E(\hat{f}_t)$ and also $\tilde{E}(\tilde{f})$. Similar to the proof of Theorem 6, we can show that

$$E(\hat{f}_t) \leq 4\sigma^2 \log(|N_\alpha|/\delta) + 2C_\alpha t + 2t\alpha\sqrt{\sigma^2 \log(2t(t+1)/\delta)},$$

holds with probability at least $1 - \delta$, where N_α is the α -covering number of \mathcal{F} .

Now we analyze $\tilde{E}(\tilde{f})$ where $\tilde{f} \in \mathcal{F}$. Similarly, a simple calculation gives that for any fixed $f \in \mathcal{F}$, $2\langle -Z, f - f_* \rangle_t$ is $2\sigma\|f - f_*\|_t$ -subgaussian. Hence, with probability $1 - \delta$, simultaneously for all $t \geq 1$,

$$\tilde{E}(f) \leq -\frac{1}{2}\|f_* - f\|_t^2 + 4\sigma^2 \log\left(\frac{1}{\delta}\right) + \frac{1}{4\sigma^2} \cdot \frac{4\sigma^2\|f - f_*\|_t^2}{2} = 4\sigma^2 \log\left(\frac{1}{\delta}\right),$$

which indicates that with probability at least $1 - \delta$, we have

$$\tilde{E}(\tilde{f}) \leq 4\sigma^2 \log\left(\frac{1}{\delta}\right).$$

Finally, as for the last term $\|\tilde{f} - f_*\|_t^2$ in (21), we have the following estimation due to the bound of the misspecification error:

$$\|\tilde{f} - f_*\|_t^2 = \sum_{p=1}^t (\tilde{f}(X_p) - f_*(X_p))^2 \leq t \cdot \zeta^2,$$

where we use the fact that $X_p = (s_p, a_p, v_p)$ satisfies that $\|v_p\|_\infty \leq H$.

We combine those bounds on the three terms in (21) above, and obtain that with probability at least $1 - 2\delta$, the following inequality holds:

$$\frac{1}{2}\|\hat{f} - f_*\|_t^2 \leq 2tC_\alpha + 2t\alpha\sqrt{\sigma^2 \log(2t(t+1)/\delta)} + 8\sigma^2 \log(2N_\alpha/\delta) + \frac{3}{2}t\zeta^2.$$

Finally, we switch δ into $\delta/2$ and multiply the above inequality by 2 on both sides. And the proof of Theorem 14 is completed. \square

Next we apply this theorem to prove the following lemma:

Lemma 15. *For any transition model P , we define its corresponding function $f_p : \mathcal{X} \rightarrow \mathbb{R}$:*

$$f_P(s, a, v) = \int P(ds'|s, a)v.$$

Then with probability at least $1 - \delta$, we have $P_* \in \mathcal{B}_k$, where

$$\mathcal{B}_k = \left\{ \tilde{P} \in \mathcal{P} : \sum_{p=1}^t (f_{\tilde{P}}(X_t) - f_{\hat{P}_t}(X_t))^2 \leq \beta_k \right\}, \quad t = k(H-1).$$

Here \hat{P}_t is defined in (9) and we choose

$$\beta_k = 8H^2 \log\left(\frac{4\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta}\right) + 4H(kH-1)\alpha \left\{ 2 + \sqrt{\log\left(\frac{8kH(kH-1)}{\delta}\right)} \right\} + 8H^3 k\epsilon^2,$$

Proof. In the following proof, the notation of X_t, Y_t, \mathcal{F} are the same as the proof of Corollary 7. We notice that $Y_t - f_P(X_t) \in [-H, H]$ for every $X_t = (s_t, a_t, V_t)$, and

$$\mathbb{E}[Y_t | \mathbb{F}_t] = \mathbb{E}[V_t(s_{t+1}) | \mathbb{F}_t] = \int P(ds' | s_t, a_t) V_t(s') = f_P(s_t, a_t, v_t) = f_P(X_t).$$

Hence $Z_t = Y_t - f_P(X_t)$ is $\frac{H}{2}$ -subgaussian given \mathbb{F}_t .

For every $f \in \mathcal{F}$, there exists some $\tilde{P} \in \mathcal{P}$ such that $f(s, a, v) = \int \tilde{P}(ds' | s, a) v(s')$, which indicates that $|f(X_t)| \leq H$. Moreover, we also have $|f_P(X_t)| \leq H$.

We next apply Theorem 14 with $C = H$ and $\sigma = \frac{H}{2}$ and $f_* = f_P$ and $\zeta = H\epsilon$ and $\tilde{f} = f_{P^*}$. According to Assumption 2, we notice that, for all $X = (s, a, v)$ with $\|v\|_\infty \leq H$, we have

$$|f_*(X) - \tilde{f}(X)| = \left| \int (P(s' | s, a) - P^*(s' | s, a)) v(s') ds' \right| \leq \|P(s' | s, a) - P^*(s' | s, a)\|_1 \|v\|_\infty \leq H\epsilon = \zeta.$$

Hence we have verified all the assumptions in Theorem 14. Hence we obtain that: for any $\alpha > 0$, with probability at least $1 - \delta$, for all $t \geq 1$, we have

$$\begin{aligned} \sum_{p=1}^t (f_P(X_p) - f_{\hat{P}_t}(X_p))^2 &\leq 4H^2 \log \left(\frac{4\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) \\ &\quad + 2H(kH - 1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{8kH(kH - 1)}{\delta} \right)} \right\} + 3H^3 k\epsilon^2. \end{aligned}$$

Moreover, noticing that

$$(f_P(X_t) - f_{P^*}(X_t))^2 = \left(\int (P(ds' | s_t, a_t) - P^*(ds' | s_t, a_t)) V_t \right)^2 \leq (H\epsilon)^2,$$

we have

$$\begin{aligned} \sum_{p=1}^t (f_{P^*}(X_p) - f_{\hat{P}_t}(X_p))^2 &\leq \sum_{p=1}^t 2(f_P(X_p) - f_{P^*}(X_p))^2 + 2(f_P(X_p) - f_{\hat{P}_t}(X_p))^2 \\ &\leq 8H^2 \log \left(\frac{4\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) + 4H(kH - 1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{8kH(kH - 1)}{\delta} \right)} \right\} + 6H^3 k\epsilon^2 + 2H^2 \epsilon^2 t \\ &\leq 8H^2 \log \left(\frac{4\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_\infty)}{\delta} \right) + 4H(kH - 1)\alpha \left\{ 2 + \sqrt{\log \left(\frac{8kH(kH - 1)}{\delta} \right)} \right\} + 8H^3 k\epsilon^2. \end{aligned}$$

which indicates that $P^* \in \mathcal{B}_k$. This finishes the proof of this corollary. \square

We next provide a lemma similar to Lemma 5, only adding the misspecification analysis.

Lemma 16. *Assuming that $P^* \in \mathcal{B}_k$, we have*

$$V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + H^2 \epsilon,$$

where

$$\xi_{h+1,k} = \langle P(\cdot | s_h^k, a_h^k), V_{h+1,k} - V_{h+1}^{\pi_k} \rangle - (V_{h+1,k}(s_{h+1}^k) - V_{h+1}^{\pi_k}(s_{h+1}^k)).$$

Note that $(\xi_{2,1}, \xi_{3,1}, \dots, \xi_{H,1}, \xi_{2,2}, \xi_{3,2}, \dots, \xi_{H,2}, \xi_{2,3}, \dots)$ is a sequence of martingale differences.

Proof. We first prove by induction that

$$V_{h,k}(s_h^k) \geq V_h^*(s_h^k) - (H+1-h)\varepsilon, \quad \forall 1 \leq h \leq H+1$$

by induction on h according to the fact that $P^* \in \mathcal{B}_k$ (but not $P \in \mathcal{B}_k$). When $h = H+1$, this inequality holds since both sides equal to 0. We assume it holds for $h+1$ and we consider the case of h . Actually we have

$$\begin{aligned} Q_{h,k}(s_h^k) &= r(s_h^k, a_h^k) + \langle P^k(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle \geq r(s_h^k, a_h^k) + \langle P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle \\ &= r(s_h^k, a_h^k) + \langle P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle - \langle P(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle \\ &\geq r(s_h^k, a_h^k) + \langle P(\cdot | s_h^k, a_h^k), V_{h+1}^* \rangle - (H-h)\varepsilon - \|P(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k)\|_1 \|V_{h+1,k}\|_\infty \\ &\geq r(s_h^k, a_h^k) + \langle P(\cdot | s_h^k, a_h^k), V_{h+1}^* \rangle - (H+1-h)\varepsilon = Q_h^*(s_h^k, a_h^k) - (H+1-h)\varepsilon, \end{aligned}$$

where in the third line we use the induction and in the last line we use the fact that $\|V_{h+1,k}\|_\infty \leq H$. This indicates that $V_{h,k}(s_h^k) \geq V_h^*(s_h^k) - (H+1-h)\varepsilon$, which completes the induction at h . Hence we know that $V_{h,k}(s_h^k) \geq V_h^*(s_h^k) - (H+1-h)\varepsilon$ holds for all $1 \leq h \leq H+1$.

Therefore,

$$V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq V_{1,k}(s_1^k) - V_1^{\pi_k}(s_1^k) + H\varepsilon.$$

Fix $h \in [H]$. In what follows we bound $V_{h,k}(s_h^k) - V_h^{\pi_k}(s_h^k)$. By the definition of π_k , P^k and a_h^k , we have

$$\begin{aligned} V_{h,k}(s_h^k) &= r(s_h^k, a_h^k) + \langle P^k(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle \text{ and} \\ V_h^{\pi_k}(s_h^k) &= r(s_h^k, a_h^k) + \langle P(\cdot | s_h^k, a_h^k), V_{h+1}^{\pi_k} \rangle. \end{aligned}$$

Hence,

$$\begin{aligned} V_{h,k}(s_h^k) - V_h^{\pi_k}(s_h^k) &= \langle P^k(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle - \langle P_{a_h^k}^{\pi_k}(s_h^k), V_{h+1}^{\pi_k} \rangle \\ &= \langle P^k(\cdot | s_h^k, a_h^k) - P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \langle P(\cdot | s_h^k, a_h^k), V_{h+1,k} - V_{h+1}^{\pi_k} \rangle. \end{aligned}$$

Therefore, by induction, noting that $V_{H+1,k} = 0$, we get that

$$\begin{aligned} V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) &\leq \sum_{h=1}^{H-1} \langle P^k(\cdot | s_h^k, a_h^k) - P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + H\varepsilon \\ &\leq \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + H\varepsilon. \end{aligned}$$

Finally noticing that

$$\begin{aligned} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle &= \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \langle \tilde{P}^*(\cdot | s_h^k, a_h^k) - P(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle \\ &\leq \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + H\varepsilon, \end{aligned}$$

we have

$$\begin{aligned} V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) &\leq \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + \frac{H(H-1)}{2} \xi + H\varepsilon \\ &\leq \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h+1,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + H^2\varepsilon, \end{aligned}$$

which completes the proof of this lemma. \square

Equipped with these two lemmas, we are ready to prove Theorem 4.

Proof of Theorem 4. According to Lemma 15, we learn that $P^* \in \mathcal{B}_k$ holds with probability at least $1 - \delta$. We next assume $P^* \in \mathcal{B}_k$ and bound the error $V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k)$.

According to Lemma 16, we have

$$V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle + \sum_{h=1}^{H-1} \xi_{h+1,k} + H^2 \varepsilon, \quad (24)$$

where

$$\xi_{h+1,k} = \langle P(\cdot | s_h^k, a_h^k), V_{h+1,k} - V_{h+1}^{\pi_k} \rangle - (V_{h+1,k}(s_{h+1}^k) - V_{h+1}^{\pi_k}(s_{h+1}^k)).$$

We let

$$W_k = \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle,$$

and summing h from 1 to H in (24) we obtain the following bound on the regret up to horizon K :

$$R_K = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \leq \sum_{k=1}^K W_k + \sum_{k=1}^K \sum_{h=1}^{H-1} \xi_{h+1,k} + H^2 K \varepsilon$$

We next bound $\sum_{k=1}^K W_k$. Actually we have

$$\begin{aligned} \sum_{k=1}^K W_k &= \sum_{k=1}^K \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle \\ &\leq \sum_{k=1}^K \sum_{h=1}^{H-1} \sup_{\tilde{P} \in \mathcal{B}_k} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle, \end{aligned}$$

and each term inside satisfies

$$\sup_{\tilde{P} \in \mathcal{B}_k} \langle \tilde{P}(\cdot | s_h^k, a_h^k) - P^*(\cdot | s_h^k, a_h^k), V_{h,k} \rangle \leq \text{diam}(\tilde{\mathcal{F}}_t | X_t)$$

where

$$\tilde{\mathcal{F}}_t = \left\{ f = f_P : P \in \mathcal{P}, \sum_{p=1}^t (f(X_p) - f_{\hat{P}_t}(X_p))^2 \leq \beta_k \right\}$$

We notice that $\mathcal{F}_t \subset \mathcal{F} \subset \mathcal{B}_\infty(\mathcal{X}, H)$. Hence we apply Lemma 9 and obtain that

$$\sum_{k=1}^K \text{diam}(\tilde{\mathcal{F}}_t | X_t) \leq \alpha + H(d \wedge K(H-1)) + 2\delta_{K(H-1)} \sqrt{dK(H-1)},$$

where $\delta_{K(H-1)} = \max_{1 \leq t \leq K(H-1)} \text{diam}(\tilde{\mathcal{F}}_t | X_t)$. Thanks to the definition of $\tilde{\mathcal{F}}_t$, $\delta_{K(H-1)} \leq 2\sqrt{\beta_K}$. Plugging this into the previous display finishes the proof.

Moreover, we also have $\sum_{k=1}^K \sum_{h=1}^{H-1} \xi_{h+1,k} \leq H\sqrt{2K(H-1)\log(1/\delta)}$ holds with probability at least $1 - \delta$. Hence combine these two inequality together, we obtain that with probability at least $1 - 2\delta$, the following bound holds

$$\begin{aligned} R_K &= \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) \\ &\leq \sum_{k=1}^K \sup_{\tilde{P} \in \mathcal{B}_k} \sum_{h=1}^{H-1} \langle \tilde{P}_{a_h^k}(s_h^k) - P_{a_h^k}(s_h^k), V_{h,k} \rangle + \sum_{k=1}^K \sum_{h=1}^{H-1} \xi_{h+1,k} + H^2 K \varepsilon \\ &\leq \alpha + H(d \wedge K(H-1)) + 4\sqrt{d\beta_K K(H-1)} + H\sqrt{2K(H-1)\log(1/\delta)} + H^2 K \varepsilon, \end{aligned}$$

where we use $\alpha \leq 1$. □

E. Proof of Theorem 13

Here we will provide the formal regret analysis for Algorithm 3, which differs from Algorithm 1. By leveraging the linear structures, we provide an independent proof of Theorem 13 using an analysis adapted from linear bandits.

The full proof is divided into five parts in the following five subsections respectively. In the first subsection, we decompose the regret into the sum of bonuses assuming the Q-functions indeed are optimistic estimates. In the second subsection, we discover some important properties of our algorithm. We provide an upper bound to the sum of bonuses in the third subsection. In the fourth subsection, we will prove that the optimism holds with high probability by constructing a particular martingale and showing that it concentrates, and in the final subsection, we will put together all the analysis to finish the proof of upper bound of expected regret.

We say $(h, k) \leq (h', k')$ if $k < k'$ or $k = k', h \leq h'$. Thus, \leq stands for the lexicographic order with k being the variable that takes priority. We say $(h, k) < (h', k')$ if $k < k'$ or $k = k', h < h'$. Let $\mathbb{F}_{h,k}$ be the filtration generated by the random sample path $\{(s_{h'}^{k'}, a_{h'}^{k'}, r_{h'}^{k'})\}_{(h', k') \leq (h, k)}$.

E.1. Regret Analysis

The proof in this section is similar to Lemma 5. Throughout E.1 to E.3, we assume that $\theta_* \in B_k$ for all $1 \leq k \leq K$. And in subsection E.4 we will prove that this event holds with high probability.

E.1.1. OPTIMISM

We will show by induction that $Q_h^*(s, a) \leq Q_{h,k}(s, a)$ for all (s, a) , h and k . For $h = H + 1$, this inequality obviously holds, since both sides equal to 0. Next suppose that this inequality holds for some $h + 1 \leq H$. As a result, we have

$$V_{h+1}^*(s) = \prod_{[0, H]} \left[\max_{a \in \mathcal{A}} Q_{h+1}^*(s, a) \right] \leq \prod_{[0, H]} \left[\max_{a \in \mathcal{A}} Q_{h+1, k}(s, a) \right] = V_{h+1, k}(s),$$

which indicates that

$$\begin{aligned} Q_h^*(s, a) &= r(s, a) + P(\cdot | s, a)^\top V_{h+1}^* \leq r(s, a) + P(\cdot | s, a)^\top V_{h+1, k} \\ &= r(s, a) + \sum_{j=1}^d (\theta_*)_j P_j(\cdot | s, a)^\top V_{h+1, k} \leq r(s, a) + \max_{\theta \in B_k} \left[\sum_{j=1}^d (\theta)_j P_j(\cdot | s, a)^\top V_{h+1, k} \right] \\ &= Q_{h, k}(s, a). \end{aligned}$$

This completes the induction.

E.1.2. REGRET DECOMPOSITION

Let us denote π_k to be the stationary policy used in the k episode, and let

$$\bar{\theta}_{h, k}(s, a) = \arg \max_{\theta \in B_k} \sum_{j=1}^d (\theta)_j P_j(\cdot | s, a)^\top V_{h+1, k}.$$

Using the fact that $\pi_k(s_h^k) = a_h^k$ and $\theta_* \in B_k$ and letting ξ_{h+1}^k be

$$\xi_{h+1}^k := P(\cdot | s_h^k, a_h^k)^\top (V_{h+1, k} - V_{h+1}^*) - [V_{h+1, k}(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k)],$$

we have

$$\begin{aligned}
 V_{h,k}(s_h^k) - V_h^{\pi_k}(s_h^k) &= Q_{h,k}(s_h^k, a_h^k) - Q_h^{\pi_k}(s_h^k, a_h^k) \\
 &= r(s_h^k, a_h^k) + \bar{\theta}_{h,k}(s_h^k, a_h^k)^\top P(\cdot | s_h^k, a_h^k) V_{h+1,k} - r(s_h^k, a_h^k) - \theta_*^\top P(\cdot | s_h^k, a_h^k) V_{h+1}^{\pi_k} \\
 &= [\theta_* + \bar{\theta}_{h,k}(s_h^k, a_h^k) - \theta_k + \theta_k - \theta_*]^\top P(\cdot | s_h^k, a_h^k) V_{h+1,k} - \theta_*^\top P(\cdot | s_h^k, a_h^k) V_{h+1}^{\pi_k} \\
 &\leq \theta_*^\top P(\cdot | s_h^k, a_h^k) (V_{h+1,k} - V_{h+1}^{\pi_k}) + 2 \max_{\theta \in B_k} |(\theta - \theta_k)^\top P(\cdot | s_h^k, a_h^k) V_{h+1,k}| \\
 &\leq P(\cdot | s_h^k, a_h^k)^\top (V_{h+1,k} - V_{h+1}^{\pi_k}) + 2 \max_{\theta \in B_k} \sqrt{(\theta - \theta_k)^\top M_k (\theta - \theta_k)} \cdot \\
 &\quad \sqrt{[P(\cdot | s_h^k, a_h^k) V_{h+1,k}]^\top M_k^{-1} [P(\cdot | s_h^k, a_h^k) V_{h+1,k}]} \\
 &\leq V_{h+1,k}(s_{h+1}^k) - V_{h+1}^{\pi_k}(s_{h+1}^k) + \xi_{h+1}^k \\
 &\quad + 2\sqrt{\beta_k} \cdot \sqrt{[P(\cdot | s_h^k, a_h^k) V_{h+1,k}]^\top M_k^{-1} [P(\cdot | s_h^k, a_h^k) V_{h+1,k}]},
 \end{aligned}$$

where the first inequality uses the fact that $\theta_*, \bar{\theta}_{h,k} \in B_k$, the second inequality uses the Cauchy-Schwarz inequality and the third inequality uses the definition of B_k .

Recall that $V_{h+1,k}(s) = V_{H+1}^*(s) = 0$ for any $s \in \mathcal{S}$. We apply the preceding inequality recursively and obtain

$$\begin{aligned}
 V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k) &\leq V_{1,k}(s_1^k) - V_1^{\pi_k}(s_1^k) \quad (\text{by optimism of value estimates}) \\
 &\leq \sum_{h=1}^H \xi_{h+1}^k + 2 \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P(\cdot | s_h^k, a_h^k) V_{h+1,k}]^\top M_k^{-1} [P(\cdot | s_h^k, a_h^k) V_{h+1,k}]},
 \end{aligned}$$

therefore the expected regret can be bounded by if we bound the expectation of

$$\begin{aligned}
 \hat{R}(K) &= \sum_{k=1}^K [V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k)] \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H \xi_{h+1}^k + 2 \sum_{k=1}^K \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P(\cdot | s_h^k, a_h^k) V_{h+1,k}]^\top M_k^{-1} [P(\cdot | s_h^k, a_h^k) V_{h+1,k}]}.
 \end{aligned} \tag{25}$$

Moreover, it is easy to observe that

$$\mathbb{E} [\xi_{h+1}^k | \mathbb{F}_{h,k}] = 0,$$

therefore ξ_{h+1}^k is a martingale difference sequence w.r.t. $\mathcal{F}_{h,k}$. Since

$$0 \leq V_h^*(s_h^k), V_{h,k}(s_h^k) \leq H \quad \text{and} \quad P(\cdot | s_h^k, a_h^k) \text{ is a probability distribution over the state space,}$$

we have $|\xi_h^k| \leq H$ with probability 1. By the Azuma-Hoeffding inequality, with probability at least $1 - \delta$, the following inequality holds

$$\sum_{k=1}^K \sum_{h=1}^H \xi_{h+1}^k \leq \sqrt{2H^3 K \log(1/\delta)}. \tag{26}$$

It remains to analyze the second term of (25), ie., the sum of bonus given by

$$2 \sum_{k=1}^K \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P(\cdot | s_h^k, a_h^k) V_{h+1,k}]^\top M_k^{-1} [P(\cdot | s_h^k, a_h^k) V_{h+1,k}]}$$

E.2. Some Properties of Algorithm 3

In this subsection we establish several useful properties of our algorithm, assuming that optimism holds throughout.

E.2.1.

Note that

$$M_{h,k} = M_{1,1} + \sum_{(h',k') < (h,k)} \left[P(\cdot | s_{h'}^{k'}, a_{h'}^{k'}) V_{h'+1,k'} \right] \left[P(\cdot | s_{h'}^{k'}, a_{h'}^{k'}) V_{h'+1,k'} \right]^\top.$$

Denote

$$l_{h,k} = \sqrt{\left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top M_{h,k}^{-1} \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]}.$$

Denote by $(h, k) + 1$ the double index of the next time step after (h, k) , that is $(h + 1, k)$ if $h < H$ and $(h, k + 1)$ otherwise. We can see $\{M_{h,k}\}$ satisfies $M_{1,k} = M_{H+1,k-1}$ and also a recursive formula

$$\begin{aligned} M_{(h,k)+1}^{-1} &= \left(M_{h,k} + \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right] \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top \right)^{-1} \\ &= M_{h,k}^{-1} - \frac{M_{h,k}^{-1} \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right] \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top M_{h,k}^{-1}}{1 + \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top M_{h,k}^{-1} \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]}. \end{aligned}$$

It implies that

$$\left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top M_{(h,k)+1}^{-1} \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right] = l_{h,k}^2 - \frac{l_{h,k}^2 \cdot l_{h,k}^2}{1 + l_{h,k}^2} = \frac{l_{h,k}^2}{1 + l_{h,k}^2}.$$

E.2.2.

Next, we derive an upper bound to the quantity

$$\sum_{k=1}^K \sum_{h=1}^H \frac{l_{h,k}^2}{1 + l_{h,k}^2}.$$

Since

$$M_{(h,k)+1} = M_{h,k} + \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right] \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top,$$

we have

$$\begin{aligned} \det M_{(h,k)+1} &= \det M_{h,k} \det \left(I + M_{h,k}^{-1/2} \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right] \left[P(\cdot | s_h^k, a_h^k) V_{h+1,k} \right]^\top M_{h,k}^{-1/2} \right) \\ &= \det M_{h,k} (1 + l_{h,k}^2). \end{aligned}$$

This indicates that

$$\sum_{(h',k') \leq (h,k)} \log(1 + l_{h',k'}^2) = \log \det M_{(h,k)+1} - \log \det M_{1,1}.$$

Furthermore, since

$$\frac{l_{h,k}^2}{1 + l_{h,k}^2} \leq \min\{1, l_{h,k}^2\} \leq 2 \log(1 + l_{h,k}^2),$$

we have

$$\begin{aligned} \sum_{(h',k') \leq (h,k)} \frac{l_{h',k'}^2}{1 + l_{h',k'}^2} &\leq \sum_{(h',k') \leq (h,k)} \min\{1, l_{h',k'}^2\} \\ &\leq \sum_{(h',k') \leq (h,k)} 2 \log(1 + l_{h',k'}^2) = 2 \log \det M_{(h,k)+1} - 2 \log \det M_{1,1}. \end{aligned}$$

E.2.3.

Given the initial value $M_{1,1} = H^2 dI$, we have

$$\begin{aligned} \mathbf{tr}(M_{(h,k)+1}) &= \mathbf{tr}(M_{1,1}) + \sum_{(h',k') \leq (h,k)} \|P_\cdot(\cdot|s_{h'}^{k'}, a_{h'}^{k'})V_{h'+1,k'}\|^2 \\ &= H^2 d^2 + \sum_{(h',k') \leq (h,k)} \sum_{j=1}^d \left(P_j(\cdot|s_{h'}^{k'}, a_{h'}^{k'})V_{h'+1,k'} \right)^2 \\ &\leq H^2 d^2 + KH^3, \end{aligned}$$

where the last inequality uses Assumption 3 and the fact that

$$P_j(\cdot|s_{h'}^{k'}, a_{h'}^{k'})V_{h'+1,k'} \leq \|P_j(\cdot|s_{h'}^{k'}, a_{h'}^{k'})\|_1 \|V_{h'+1,k'}\|_\infty \leq H.$$

Using the inequalities of arithmetic and geometric means, we get the following upper bound for the determinant of $M_{(h,k)+1}$:

$$\det M_{(h,k)+1} \leq \left(\frac{\mathbf{tr}(M_{(h,k)+1})}{d} \right)^d \leq (H^2 d + KH^3)^d,$$

which indicates that

$$\log \det M_{(H,k)+1} - \log \det M_{1,1} \leq \log((H^2 d + KH^3)^d) - \log((H^2 d)^d) \leq d \log(1 + HK). \quad (27)$$

Hence we have

$$\sum_{(h',k') \leq (h,k)} \frac{l_{h',k'}^2}{1 + l_{h',k'}^2} \leq \sum_{(h',k') \leq (h,k)} \min\{1, l_{h',k'}^2\} \leq 2d \log(1 + HK). \quad (28)$$

E.3. Sum-of-Bonus Analysis

In this section, under the assumption that $\theta_* \in B_k$ for every k , we establish an upper bound for the following sum-of-bonus term

$$2 \sum_{k=1}^K \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]},$$

where we denote $M_k = M_{1,k}$ for simplicity. We let

$$u_{h,k} = \sqrt{[P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]}.$$

Since $\beta_k \leq \beta_K$ for any $1 \leq k \leq K$ and by letting

$$\begin{aligned} u_{h,k}^2 &= [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}] \\ &\leq [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_1^{-1} [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}] \\ &= \frac{1}{H^2 d} \cdot \sum_{j=1}^d [P_j(\cdot|s_h^k, a_h^k)V_{h+1,k}]^2 \leq \frac{1}{H^2 d} \cdot H^2 d = 1, \end{aligned}$$

we have

$$\begin{aligned} &2 \sum_{k=1}^K \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P_\cdot(\cdot|s_h^k, a_h^k)V_{h+1,k}]} \\ &\leq 2\sqrt{\beta_K} \cdot \sum_{k=1}^K \sum_{h=1}^H u_{h,k} \leq 2\sqrt{\beta_K} \cdot \sum_{k=1}^K \sum_{h=1}^H \min\{1, u_{h,k}\} \\ &\leq 2\sqrt{HK\beta_K} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \min\{1, u_{h,k}^2\}} \leq 4\sqrt{HK\beta_K} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \log(1 + u_{h,k}^2)} \end{aligned} \quad (29)$$

where the third inequality uses the Cauchy-Schwarz inequality. Next we notice that

$$M_{k+1} = M_k + \sum_{h=1}^H [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}].$$

Hence we have

$$\det(M_{k+1}) = \det(M_k) \cdot \det\left(I + \sum_{h=1}^H M_k^{-1/2} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] M_k^{-1/2}\right).$$

We further notice that every eigenvalue of the matrix

$$I + \sum_{h=1}^H M_k^{-1/2} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] M_k^{-1/2}$$

is at least 1, and we have the following bound of its trace:

$$\begin{aligned} & \text{tr}\left(\sum_{h=1}^H M_k^{-1/2} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] M_k^{-1/2}\right) \\ &= \sum_{h=1}^H [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] = \sum_{h=1}^H u_{h,k}^2. \end{aligned}$$

This indicates that

$$\begin{aligned} & \det\left(I + \sum_{h=1}^H M_k^{-1/2} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] M_k^{-1/2}\right) \\ & \geq 1 + \text{tr}\left(I + \sum_{h=1}^H M_k^{-1/2} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top [P(\cdot|s_h^k, a_h^k)V_{h+1,k}] M_k^{-1/2}\right) \\ & = 1 + \sum_{h=1}^H u_{h,k}^2, \end{aligned}$$

where the first inequality follows from the following fact: $\prod_i(1+w_i) \geq 1 + \sum_i w_i$ provided $w_i \geq 0$. Combining the above inequality with the following inequality

$$1 + \sum_{h=1}^H u_{h,k}^2 = \frac{\sum_{h=1}^H (1 + H u_{h,k}^2)}{H} \geq \prod_{h=1}^H (1 + H u_{h,k}^2)^{1/H} \geq \prod_{h=1}^H (1 + u_{h,k}^2)^{1/H},$$

we obtain that

$$\sum_{h=1}^H \log(1 + u_{h,k}^2) \leq H \log\left(1 + \sum_{h=1}^H u_{h,k}^2\right) \leq H \det(M_{k+1}) - H \det(M_k).$$

Therefore, we have

$$\begin{aligned} & 2 \sum_{k=1}^K \sum_{h=1}^H \sqrt{\beta_k} \cdot \sqrt{[P(\cdot|s_h^k, a_h^k)V_{h+1,k}]^\top M_k^{-1} [P(\cdot|s_h^k, a_h^k)V_{h+1,k}]} \\ & \leq 4\sqrt{HK\beta_K} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \log(1 + u_{h,h}^2)} \\ & \leq 4\sqrt{HK\beta_K} \cdot \sqrt{\sum_{k=1}^K H \det(M_{k+1}) - H \det(M_k)} \\ & \leq 4\sqrt{HK\beta_K} \cdot \sqrt{H \det(M_{(H,k)+1}) - H \det(M_{1,1})} \\ & \leq 4\sqrt{H^2 d K \beta_K \log(1 + HK)}, \end{aligned}$$

where the last inequality uses (27).

E.4. Confidence Set of Value Target Regression

We adopt the result in (Abbasi-Yadkori et al., 2011). For $t = H(k-1) + h$, we choose

$$\begin{aligned}\lambda &= H^2 d, \\ \bar{V}_t &= M_{h,k}, \\ S &= C_\theta, \\ R &= H, \\ L &= \sqrt{H^2 d}.\end{aligned}$$

Then we have

$$\theta_{h,k} = (\mathbf{X}_{1:t}^T \mathbf{X}_{1:t} + \lambda I)^{-1} \mathbf{X}_{1:t} \mathbf{Y}_{1:t} = \hat{\theta}_t, \quad \text{and} \quad \|\theta_*\|_2 \leq C_\theta = S.$$

Moreover, if we since $|\eta_t| = |Y_t - \langle X_t, \theta_* \rangle| = |Y_{h,k} - P(\cdot | s_h^k, a_h^k)^T V_{h+1,k}| \leq H$, η_t is H -subgaussian. We can also verify that

$$\|X_t\|_2^2 = \sum_{i=1}^d (P_i(\cdot | s_h^k, a_h^k) V_{h+1,k})^2 \leq H^2 d = L^2.$$

Hence according to Theorem 2 in (Abbasi-Yadkori et al., 2011), we obtain that with probability at least $1 - \delta$, for any $(h, k) \leq (H, K)$, the following inequality holds:

$$\|\theta_* - \theta_{h,k}\|_{M_{h,k}} \leq H \sqrt{d \log \left(\frac{1 + Hk \cdot H^2 d}{\delta} \right)} + C_\theta H \sqrt{d}$$

Therefore, if we choose

$$\beta_k = \left(H \sqrt{d \log \left(\frac{1 + Hk \cdot H^2 d}{\delta} \right)} + C_\theta H \sqrt{d} \right)^2,$$

then we will have

$$\theta_* \in B_{h,k}$$

for all $(h, k) \leq (H, K)$ with probability at least $1 - \delta$.

E.5. Expected Regret Analysis

According to Section E.4, we have with probability at least $1 - \delta$ that $\theta_* \in B_k$ for all $1 \leq k \leq K$. When this event happens, we enable the analysis of Sections E.1-E.3. We combine the error bounds (26) and (29) and apply them into the regret bound (25). It follows that, if $T = KN$,

$$\begin{aligned}R(T) &\leq 2\sqrt{H^3 K} \log(1/\delta) + 4\sqrt{H^2 d K \beta_K \log(1 + HK)} \\ &= 2\sqrt{H^3 K} \log\left(\frac{1}{\delta}\right) + 4H^2 d \sqrt{K \log(1 + HK)} \cdot \left(\sqrt{\log\left(\frac{1 + H^3 K d}{\delta}\right)} + C_\theta \right) \\ &\leq 6H^2 d \sqrt{K} \left(C_\theta \sqrt{\log(1 + HK)} + \log\left(\frac{1 + H^3 K d}{\delta}\right) \right)\end{aligned}$$

with probability at least $1 - 2\delta$. Note the trivial upper bound $R(K) \leq HK$. Therefore, by letting $\delta = 1/K$ and noticing $T = HK$, we get

$$\begin{aligned}\mathbb{E}[R(T)] &\leq (1 - 2\delta) \cdot 6H^2 d \sqrt{K} \left(C_\theta \sqrt{\log(1 + HK)} + \log\left(\frac{1 + H^3 K d}{\delta}\right) \right) + 2\delta \cdot HK \\ &\leq 6H^2 d \sqrt{K} \cdot \left(C_\theta \sqrt{\log(1 + HK)} + \log\left(\frac{1 + H^3 K d}{\delta}\right) \right) \\ &= \tilde{O}(C_\theta \cdot H^2 d \sqrt{K}) = \tilde{O}(C_\theta \cdot d \sqrt{H^3 T}).\end{aligned}$$

Thus we have completed the proof of Theorem 13.

F. Implementation

F.1. Analysis of Implemented Confidence Bounds

In the implementation of UCRL-VTR used in Section 6, we used different confidence intervals than the ones stated in the paper. The confidence intervals used in our implementation are the ones introduced in (Abbasi-Yadkori et al., 2011). These confidence intervals are much tighter in the linear setting than the ones introduced in Section 3 and thus have better practical performance. The purpose of this section is to formally introduce the confidence intervals used in our implementation of UCRL-VTR as well as show how these confidence intervals were adapted from the linear bandit setting to the linear MDP setting.

F.1.1. LINEAR MDP ASSUMPTIONS

For our implementation of UCRL-VTR we used different confidence than was introduced in the paper. These are the tighter confidence bounds from the seminal work done by (Abbasi-Yadkori et al., 2011) and further expanded upon in Chapter 20 of (Lattimore & Szepesvári, 2020). Now we will state some assumptions in the MDP setting, then we will state the equivalent assumptions from the linear bandit setting, and lastly we will make the connections between the two that allow us to use the confidence bounds from the linear bandit setting in the RL setting.

1. $P^*(s' | s, a) = \sum_{i=1}^d (\theta_*^{MDP})_i P_i(s' | s, a)$
2. $s_{h+1}^k \sim P^*(\cdot | s_h^k, a_h^k)$
3. $\mathcal{C}_t^{MDP} = \{\theta^{MDP} \in \mathbb{R}^d : \|\theta^{MDP} - \hat{\theta}_t^{MDP}\|_{M_t} \leq \beta_t\}$

where t is defined in the table of A.4. Also note that in this section $(\cdot)_*$ denotes the true parameter or model, $(\cdot)^{MDP}$ denotes something derived or used in the linear MDP setting, and $(\cdot)^{LIN}$ denotes something derived or used in the linear bandit setting. Now, under 1-3 of F.1.1 we hope to construct a confidence set \mathcal{C}_t^{MDP} such that

$$\theta^{MDP} \in \bigcap_{t=1}^{\infty} \mathcal{C}_t^{MDP}$$

with high probability. Now the choice of how to choose both \mathcal{C}_t^{MDP} and β_t comes from the linear bandit literature. We will introduce the necessary theorems and assumptions to derive both \mathcal{C}_t^{LIN} and β_t in the linear bandit setting and then adapt the results from the linear bandit setting to the linear MDP setting.

F.1.2. TIGHTER CONFIDENCE BOUNDS FOR LINEAR BANDITS

The following results are introduced in the paper by (Abbasi-Yadkori et al., 2011) and are further explained in Chapter 20 of the book by (Lattimore & Szepesvári, 2020). In this section, we will introduce the theorems and lemmas that allows us to derive tighter confidence intervals for the linear bandit setting. Then we will carefully adapt the confidence intervals to the linear bandit setting. Now supposed a bandit algorithm has chosen actions $A_1, \dots, A_t \in \mathbb{R}^d$ and received rewards $X_1^{LIN}, \dots, X_t^{LIN}$ with $X_s^{LIN} = \langle A_s, \theta_*^{LIN} \rangle + \eta_s$ where η_s is some zero mean noise. The least squares estimator of θ_*^{LIN} is the minimizer of the following loss function

$$L_t(\theta^{LIN}) = \sum_{s=1}^t (X_s^{LIN} - \langle A_s, \theta^{LIN} \rangle)^2 + \lambda \|\theta^{LIN}\|_2^2$$

where $\lambda > 0$ is the regularizer. This loss function is minimized by

$$\hat{\theta}_t^{LIN} = W_t^{-1} \sum_{s=1}^t X_s^{LIN} A_s \text{ with } W_t = \lambda I + \sum_{s=1}^t A_s A_s^\top$$

notice how this linear bandit problem is very similar to the linear MDP problem introduced in section 3 of our paper. In our linear MDP setting, it is convenient to think of M and W as serving equivalent purposes (storing rank one updates) thus it is also convenient to think of A_t and X_t^{MDP} as serving equivalent purposes (the features by which we use to make our predictions), where X_t^{MDP} is defined in section 3 of our paper with some added notation to distinguish it from the X_t^{LIN} used here in the linear bandit setting. We will now build up some intuition by making some simplifying assumptions.

1. No regularization: $\lambda = 0$ and W_t is invertible.
2. Independent subgaussian noise: $(\eta_s)_s$ are independent and σ -subgaussian
3. Fixed Design: A_1, \dots, A_t are deterministically chosen without the knowledge of $X_1^{LIN}, \dots, X_t^{LIN}$

finally it is also convenient to think of X_t^{LIN} and $V_{t+1}(s_{t+1})$ as serving equivalent purposes (the target of our predictions). Thus the statements we prove in the linear bandit setting can be easily adapted to the linear MDP setting. While none of the assumptions stated above is plausible in the bandit setting, the simplifications eases the analysis and provides insight.

Comparing θ_*^{LIN} and $\hat{\theta}_t^{LIN}$ in the direction $x \in \mathbb{R}^d$, we have

$$\begin{aligned} \langle \hat{\theta}_t^{LIN} - \theta_*^{LIN}, x \rangle &= \left\langle x, W_t^{-1} \sum_{s=1}^t A_s X_s^{LIN} - \theta_*^{LIN} \right\rangle = \left\langle x, W_t^{-1} \sum_{s=1}^t A_s (A_s^\top \theta_*^{LIN} + \eta_s) - \theta_*^{LIN} \right\rangle \\ &= \left\langle x, W_t^{-1} \sum_{s=1}^t A_s \eta_s \right\rangle = \sum_{s=1}^t \langle x, W_t^{-1} A_s \rangle \eta_s \end{aligned}$$

Since $(\eta_s)_s$ are independent and σ -subgaussian, by Lemma 5.4 and Theorem 5.3 (need to be stated),

$$\mathbb{P} \left(\langle \hat{\theta}_t^{LIN} - \theta_*^{LIN}, x \rangle \geq \sqrt{2\sigma^2 \sum_{s=1}^t \langle x, W_t^{-1} A_s \rangle^2 \log \left(\frac{1}{\delta} \right)} \right) \leq \delta$$

A little linear algebra shows that $\sum_{s=1}^t \langle x, W_t^{-1} A_s \rangle^2 = \|x\|_{W_t^{-1}}^2$ and so,

$$\mathbb{P} \left(\langle \hat{\theta}_t^{LIN} - \theta_*^{LIN}, x \rangle \geq \sqrt{2\sigma^2 \|x\|_{W_t^{-1}}^2 \log \left(\frac{1}{\delta} \right)} \right) \leq \delta \quad (30)$$

We now remove the limiting assumptions we stated above and use the newly stated assumptions for the rest of this section

1. There exists a $\theta_*^{LIN} \in \mathbb{R}^d$ such that $X_t^{LIN} = \langle \theta_*^{LIN}, A_t \rangle + \eta_t$ for all $t \geq 1$.
2. The noise is conditionally σ -subgaussian:

$$\text{for all } \alpha \in \mathbb{R} \text{ and } t \geq 1, \mathbb{E}[\exp(\alpha \eta_t) \mid \mathcal{F}_{t-1}] \leq \exp \left(\frac{\alpha \sigma^2}{2} \right) a.s.$$

where \mathcal{F}_{t-1} is such that $A_1, X_1^{LIN}, \dots, A_{t-1}, X_{t-1}^{LIN}$ are \mathcal{F}_{t-1} -measurable.

3. In addition, we now assume $\lambda > 0$.

The inclusion of A_t in the definition of \mathcal{F}_{t-1} allows the noise to depend on past choices, including the most recent action. Since we want exponentially decaying tail probabilities, one is tempted to try the Cramer-Chernoff method:

$$\mathbb{P}(\|\hat{\theta}_t^{LIN} - \theta_*^{LIN}\|_{W_t}^2 \geq u^2) \leq \inf_{\alpha > 0} \mathbb{E} \left[\exp \left(\alpha \|\hat{\theta}_t^{LIN} - \theta_*^{LIN}\|_{W_t}^2 - \alpha u^2 \right) \right].$$

Sadly, we do not know how to bound this expectation. Can we still somehow use the Cramer-Chernoff method? We take inspiration from looking at the special case of $\lambda = 0$ one last time, assuming that $W_t = \sum_{s=1}^t A_s A_s^\top$ is invertible. Let

$$S_t = \sum_{s=1}^t \eta_s A_s$$

Recall that $\hat{\theta}_t^{LIN} = W_t^{-1} \sum_{s=1}^t X_s^{LIN} A_s = \theta_*^{LIN} + W_t^{-1} S_t$. Hence,

$$\frac{1}{2} \|\hat{\theta}_t^{LIN} - \theta_*^{LIN}\|_{W_t}^2 = \frac{1}{2} \|S_t\|_{W_t^{-1}}^2 = \max_{x \in \mathbb{R}^d} \left(\langle x, S_t \rangle - \frac{1}{2} \|x\|_{W_t}^2 \right).$$

The next lemma shows that the exponential of the term inside the maximum is a supermartingale even when $\lambda \geq 0$.

Lemma 17. For all $x \in \mathbb{R}^d$ the process $D_t(x) = \exp(\langle x, S_t \rangle - \frac{1}{2}\|x\|_{W_t}^2)$ is an \mathbb{F} -adapted non-negative supermartingale with $D_0(x) \leq 1$.

The proof for this Lemma can be found in Chapter 20 of the book by (Lattimore & Szepesvári, 2020). For simplicity, consider now again the case when $\lambda = 0$. Combining the lemma and the linearisation idea almost works. The Cramer–Chernoff method leads to

$$\mathbb{P}\left(\frac{1}{2}\|\hat{\theta}_t^{LIN} - \theta_*^{LIN}\|_{W_t}^2 \geq \log(1/\delta)\right) = \mathbb{P}\left(\exp\left(\max_{x \in \mathbb{R}^d} \left(\langle x, S_t \rangle - \frac{1}{2}\|x\|_{W_t}^2\right)\right) \geq \log(1/\delta)\right) \quad (31)$$

$$\leq \delta \mathbb{E}\left[\exp\left(\max_{x \in \mathbb{R}^d} \left(\langle x, S_t \rangle - \frac{1}{2}\|x\|_{W_t}^2\right)\right)\right] = \delta \mathbb{E}\left[\max_{x \in \mathbb{R}^d} D_t(x)\right] \quad (32)$$

Now Lemma 17 shows that $\mathbb{E}[D_t(x)] \leq 1$. Now using Laplace’s approximation we write

$$\max_x D_t(x) \approx \int_{\mathbb{R}^d} D_t(x) dh(x),$$

where h is some measure on \mathbb{R}^d chosen so that the integral can be calculated in closed form. This is not a requirement of the method, but it does make the argument shorter. The main benefit of replacing the maximum with an integral is that we obtain the following lemma

Lemma 18. Let h be a probability measure on \mathbb{R}^d ; then; $\bar{D}_t = \int_{\mathbb{R}^d} D_t(x) dh(x)$ is an \mathbb{F} -adapted non-negative supermartingale with $\bar{D}_0 = 1$.

The proof of Lemma 18 can, again, be found in Chapter 20 of the book by (Lattimore & Szepesvári, 2020). Now the following theorem is the key result from which the confidence set will be derived.

Theorem 19. For all $\lambda > 0$, and $\delta \in (0, 1)$

$$\mathbb{P}\left(\text{exists } t \in \mathbb{N} : \|S_t\|_{W_t^{-1}}^2 \geq 2\sigma^2 \log\left(\frac{1}{\delta}\right) + \log\left(\frac{\det W_t}{\lambda^d}\right)\right) \leq \delta$$

Furthermore, if $\|\theta_*^{LIN}\|_2 \leq m_2$, then $\mathbb{P}(\text{exists } t \in \mathbb{N}^+ : \theta_*^{LIN} \notin \mathcal{C}_t^{LIN}) \leq \delta$ with

$$\mathcal{C}_t^{LIN} = \left\{ \theta \in \mathbb{R}^d : \|\hat{\theta}_{t-1}^{LIN} - \theta\|_{W_{t-1}} < m_2 \sqrt{\lambda} + \sqrt{2\sigma^2 \log\left(\frac{1}{\delta}\right) + \log\left(\frac{W_{t-1}}{\lambda^d}\right)} \right\}.$$

The proof of Theorem 19 can be found in Chapter 20 of the book by (Lattimore & Szepesvári, 2020).

F.1.3. ADAPTATION OF THE CONFIDENCE BOUNDS TO OUR LINEAR MDP SETTING

Now with the Lemmas and Theorems introduced in the previous section we are ready to derive the confidence bounds used in our implementation of UCRL-VTR. Now using the notation from the linear bandit setting we set

1. The target $X_t^{MDP} = \int_j V_t(s') P_j(ds' | s_t, a_t)$
2. $Y_t = V_t(s_{t+1})$
3. $\mathcal{F}_{t-1} = \sigma(s_1, a_1, \dots, s_{t-1}, a_{t-1})$, which just means the filtration is set to be the sigma-algebra generated by all past states and actions observed.
4. $\eta_t = Y_t - \langle X_t^{MDP}, \theta_*^{MDP} \rangle = V_t(s_{t+1}) - \int_j V_t(s') P_j^*(ds' | s_t, a_t)$, since θ_*^{MDP} is the true model of the MDP.
5. M_t in the linear MDP setting is defined equivalently to W_t in the linear bandit setting, i.e. they are both the sums of a regularizer term and a bunch of rank one updates.

it can be seen that our the noise in our system η_t has zero mean $\mathbb{E}[\eta_t | \mathcal{F}_{t-1}] = 0$ finally the noise in our system has variance $H/2$ thus our system in $H/2$ -subgaussian.

Lemma 20. (Hoeffding's lemma) Let $Z = Z - \mathbb{E}[Z]$ be a real centered random variable such that $Z \in [a, b]$ almost surely. Then $\mathbb{E}[\exp(\alpha Z)] \leq \exp(\alpha^2 \frac{(b-a)^2}{8})$ for any $\alpha \in \mathbb{R}$ or Z is subgaussian with variance $\sigma^2 = \frac{(b-a)^2}{4}$.

Proof Define $\psi(\alpha) = \log \mathbb{E}[\exp(\alpha Z)]$ we can then compute

$$\psi'(\alpha) = \frac{\mathbb{E}[Z \exp(\alpha Z)]}{\mathbb{E}[\exp(\alpha Z)]}, \quad \psi''(\alpha) = \frac{\mathbb{E}[Z^2 \exp(\alpha Z)]}{\mathbb{E}[\exp(\alpha Z)]} - \left(\frac{\mathbb{E}[Z \exp(\alpha Z)]}{\mathbb{E}[\exp(\alpha Z)]} \right)^2$$

Thus $\psi''(\alpha)$ can be interpreted as the variance of the random variable Z under the probability measure $d\mathbb{Q} = \frac{\exp(\alpha Z)}{\mathbb{E}[\exp(\alpha Z)]} d\mathbb{P}$, but since $Z \in [a, b]$ almost surely, we have, under any probability

$$\text{var}(Z) = \text{var}\left(Z - \frac{a+b}{2}\right) \leq \mathbb{E}\left[\left(Z - \frac{a+b}{2}\right)^2\right] \leq \left(\frac{b-a}{4}\right)^2$$

The fundamental theorem of calculus yields

$$\psi(\alpha) = \int_0^\alpha \int_0^\mu \psi(\rho) d\rho d\mu = \frac{s^2(b-a)^2}{8}$$

using $\psi(0) = \log 1 = 0$ and $\psi'(0) = \mathbb{E}[Z] = 0$. □

Now using Lemma 20 and the fact that Y_t is bounded in the range of $[0, H]$, $\mathbb{E}[Y_t] = \langle X_t^{MDP}, \theta_*^{MDP} \rangle$, and $\eta_t = Y_t - \langle X_t^{MDP}, \theta_*^{MDP} \rangle = Y_t - \mathbb{E}[Y_t]$, the noise η_t in our linear MDP setting is $H/2$ -subgaussian. This result is also stated in a proof from A.4.

Putting this all together we can derive the tighter confidence set for UCRL-VTR in the linear setting,

$$\mathcal{C}_t^{MDP} = \left\{ \theta \in \mathbb{R}^d : \|\hat{\theta}_{t-1}^{MDP} - \theta\|_{M_{t-1}} < m_2 \sqrt{\lambda} + \frac{H}{2} \sqrt{2 \log\left(\frac{1}{\delta}\right) + \log\left(\frac{M_{t-1}}{\lambda^d}\right)} \right\}.$$

where here in the linear MDP setting M_t replaces W_t from the linear bandit setting and $\|\theta_*^{MDP}\|_2 \leq m_2$. The justification of using these bounds in the linear MDP setting follow exactly from the justification given above for using these bounds in the linear bandit setting.

F.2. UCRL-VTR

In the proceeding subsections we discuss the implementation of the algorithms studied in Section 6 of the paper. The first algorithm we present is the algorithm used to generate the results for UCRL-VTR.

Algorithm 4 UCRL-VTR with Tighter Confidence Bounds

- 1: **Input:** MDP, $d, H, T = KH$;
- 2: **Initialize:** $M_{1,1} \leftarrow I$, $w_{1,1} \leftarrow 0 \in \mathbb{R}^{d \times 1}$, $\theta_1 \leftarrow M_{1,1}^{-1} w_{1,1}$ for $1 \leq h \leq H$, $d_1 = |\mathcal{S}| \times |\mathcal{A}|$;
- 3: **Initialize:** $\delta \leftarrow 1/K$, and for $1 \leq k \leq K$,
- 4: Compute Q-function $Q_{h,1}$ using $\theta_{1,1}$ according to (3);
- 5: **for** $k = 1 : K$ **do**
- 6: Obtain initial state s_1^k for episode k ;
- 7: **for** $h = 1 : H$ **do**
- 8: Choose action greedily by

$$a_h^k = \arg \max_{a \in \mathcal{A}} Q_{h,k}(s_h^k, a)$$
 and observe the next state s_{h+1}^k .
- 9: Compute the predicted value vector: ▷ Evaluate the expected value of next state

$$X_{h,k} \leftarrow \mathbb{E}.[V_{h+1,k}(s) | s_h^k, a_h^k] = \sum_{s \in \mathcal{S}} V_{h+1,k}(s) \cdot P.(s | s_h^k, a_h^k).$$
- 10: $y_{h,k} \leftarrow V_{h+1,k}(s_{h+1}^k)$ ▷ Update regression parameters
- 11: $M_{h+1,k} \leftarrow M_{h,k} + X_{h,k} X_{h,k}^\top$
- 12: $w_{h+1,k} \leftarrow w_{h,k} + y_{h,k} \cdot X_{h,k}$
- 13: **end for**
- 14: Update at the end of episode: ▷ Update Model Parameters

$$\begin{aligned} M_{1,k+1} &\leftarrow M_{H+1,k}, \\ w_{1,k+1} &\leftarrow w_{H+1,k}, \\ \theta_{k+1} &\leftarrow M_{1,k+1}^{-1} w_{1,k+1}; \end{aligned}$$
- 15: Compute $Q_{h,k+1}$ for $h = H, \dots, 1$, using θ_{k+1} according to (33) using

$$\sqrt{\beta_{h,k}} \leftarrow \sqrt{d_1} + \frac{H-h+1}{2} \sqrt{2 \log \left(\frac{1}{\delta} \right) + \log \det(M_{1,k+1})};$$
▷ Computing Q functions
- 16: **end for**

The iterative Q-update for Algorithm 4 is

$$\begin{aligned} V_{h+1,k}(s) &= 0 \\ Q_{h,k}(s, a) &= r(s, a) + X_{h,k}^\top \theta_k + \sqrt{\beta_{h,k}} \sqrt{X_{h,k}^\top M_{1,k+1}^{-1} X_{h,k}} \\ V_{h,k}(s) &= \max_a Q_{h,k}(s, a) \end{aligned} \tag{33}$$

The choice of the confidence bounds used in Algorithm 4 comes from the tight bounds derived in (Abbasi-Yadkori et al., 2011) for linear bandits and further expanded upon in Chapter 20 of (Lattimore & Szepesvári, 2020). The details of which are shown and stated in F.1. We slightly tighten the values for the noise at each stage by using the fact that for each stage in the horizon, $h \in [H]$, the value $V_h^k(\cdot)$ is capped as to never be greater than $H - h + 1$. The appearance of the $\sqrt{d_1}$ comes from the fact that $\|\theta_*\|_2 \leq \sqrt{d_1}$ for all $\theta_* \in \mathbb{R}^d$ in the tabular setting since θ_* in the tabular setting is equal to the true model of the environment.

F.3. EGRL-VTR

In this section we discuss the algorithm EGRL-VTR. This algorithm is very similar to UCRL-VTR except it performs ε -greedy value iteration instead of optimistic value iteration and acts ε -greedy with respect to $Q_{h,k}$.

Algorithm 5 EGRL-VTR

-
- 1: **Input:** MDP, $d, H, T = KH, \varepsilon > 0$;
 2: **Initialize:** $M_{1,1} \leftarrow I, w_{1,1} \leftarrow 0 \in \mathbb{R}^{d \times 1}, \theta_1 \leftarrow M_{1,1}^{-1} w_{1,1}$ for $1 \leq h \leq H$;
 3: Compute Q-function $Q_{h,1}$ using $\theta_{1,1}$ according to (34);
 4: **for** $k = 1 : K$ **do**
 5: Obtain initial state s_1^k for episode k ;
 6: **for** $h = 1 : H$ **do**
 7: With probability $1 - \varepsilon$ do

$$a_h^k = \arg \max_{a \in \mathcal{A}} Q_{h,k}(s_h^k, a)$$

 else pick a uniform random action $a_h^k \in \mathcal{A}$. Observe the next state s_{h+1}^k .
 8: Compute the predicted value vector: ▷ Evaluate the expected value of next state

$$X_{h,k} \leftarrow \mathbb{E}.[V_{h+1,k}(s) | s_h^k, a_h^k] = \sum_{s \in \mathcal{S}} V_{h+1,k}(s) \cdot P.(s | s_h^k, a_h^k).$$

 9: $y_{h,k} \leftarrow V_{h+1,k}(s_{h+1}^k)$ ▷ Update regression parameters
 10: $M_{h+1,k} \leftarrow M_{h,k} + X_{h,k} X_{h,k}^\top$
 11: $w_{h+1,k} \leftarrow w_{h,k} + y_{h,k} \cdot X_{h,k}$
 12: **end for**
 13: Update at the end of episode: ▷ Update Model Parameters

$$M_{1,k+1} \leftarrow M_{H+1,k},$$

$$w_{1,k+1} \leftarrow w_{H+1,k},$$

$$\theta_{k+1} \leftarrow M_{1,k+1}^{-1} w_{1,k+1};$$

 14: Compute $Q_{h,k+1}$ for $h = H, \dots, 1$, using θ_{k+1} according to (34) ▷ Computing Q functions
 15: **end for**
-

The iterative value update for EGRL-VTR is

$$\begin{aligned}
 V_{h+1,k}(s) &= 0 \\
 Q_{h,k}(s, a) &= r(s, a) + X_{h,k}^\top \theta_k \\
 V_{h,k}(s) &= (1 - \varepsilon) \Pi_{[0,H]} \max_a Q_{h,k}(s, a) + \frac{\varepsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} Q_{h,k}(s, a)
 \end{aligned} \tag{34}$$

E.4. EG-Frequency

In this section we discuss the algorithm EG-Frequency. This algorithm is the ε -greedy version of UC-MatrixRL (Yang & Wang, 2019a).

Algorithm 6 EG-Frequency

```

1: Input: MDP, Features  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  and  $\psi : \mathcal{S} \rightarrow \mathbb{R}^{|\mathcal{S}|}$ ,  $\varepsilon > 0$ , and the total number of episodes  $K$ ;
2: Initialize:  $A_1 \leftarrow I \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ ,  $M_1 \leftarrow 0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ , and  $K_\psi \leftarrow \sum_{s' \in \mathcal{S}} \psi(s')\psi(s')^\top$ ;
3: for  $k = 1 : K$  do
4:   Let  $Q_{h,k}$  be given in (35) using  $M_k$ ;
5:   for  $h = 1 : H$  do
6:     Let the current state be  $s_h^k$ ;
7:     With probability  $(1-\varepsilon)$  play action  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_{h,k}(s_h^k, a)$  else pick a uniform random action  $a_h^k \in \mathcal{A}$ .
8:     Record the next state  $s_{h+1}^k$ 
9:   end for
10:   $A_{k+1} \leftarrow A_k + \sum_{h \leq H} \phi(s_h^k, a_h^k)\phi(s_h^k, a_h^k)^\top$ 
11:   $M_{k+1} \leftarrow M_k + A_{k+1}^{-1} \sum_{h \leq H} \phi(s_h^k, a_h^k)\psi(s_{h+1}^k)^\top K_\psi^{-1}$ 
12: end for

```

The iterative Q-update for EG-Frequency is

$$\begin{aligned}
 Q_{h+1,k}(s, a) &= 0 \text{ and} \\
 Q_{h,k}(s, a) &= r(s, a) + \phi(s, a)^\top M_k \Psi^\top V_{h+1,k} \\
 V_{h,k} &= (1 - \varepsilon) \Pi_{[0, H]} \max_a Q_{h,k}(s, a) + \frac{\varepsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} Q_{h,k}(s, a)
 \end{aligned} \tag{35}$$

Note that Ψ is a $|\mathcal{S}| \times |\mathcal{S}|$ whose rows are the features $\psi(s')$ and Φ is a $|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|$ whose rows are the features $\phi(s, a)$. In the tabular RL setting both Ψ and Φ are the identity matrix which is what we used in our numerical experiments. In the tabular RL setting, EG-Frequency stores the counts of the number of times it transitioned to next state s' from the state-action pair (s, a) and fits the estimated model M_k accordingly.

F.5. Futher Implementation Notes

In this section, we include some further details on how we implemented Algorithms 4, 5, and 6. All code was written in Python 3 and used the Numpy and Scipy libraries. All plots were generated using Matplotlib. In Algorithm 4, Numpy's logdet function was used to calculate the determinate in step 15 for numerical stability purposes. No matrix inversion was performed in our code, instead a Sherman-Morrison update was performed for each matrix in which a matrix inversion is performed at each (k, h) in order to save on computation. To read more about the Sherman Morrison update in the context of RL, we refer to the reader to Eqn (9.22) of (Sutton & Barto, 2018). When computing the weighted L1-norm, we added a small constant to each summation in the denominator to avoid dividing by zero. Finally, when computing UC-MatrixRL we also used the self-normalize bounds introduced in the beginning of this section. Some pseudocode for using self-normalized bounds with UC-MatrixRL can be found in step 5 of Alg 7.

G. Mixture Model

In this section, we introduce, analyze, and evaluate a Linear model-based RL algorithm that used both the canonical model and the VTR model for planning. We call this algorithm UCRL-MIX.

G.1. UCRL-MIX

Below a meta-algorithm for UCRL-MIX

Algorithm 7 UCRL-MIX

- 1: Compute Algorithm 4 and UC-MatrixRL (Yang & Wang, 2019a) simultaneously.
 - 2: At end of episode k , perform value iteration and set $V_{H+1,k}(s) = 0$.
 - 3: **for** $h = H : 1$ **do**
 - 4: **for** $s \in |\mathcal{S}|$ and $a \in |\mathcal{A}|$ **do**
 - 5: Compute the confidence set bonuses as follows

$$B_{h,k}^{VTR} \leftarrow \sqrt{d_1} + \frac{H-h+1}{2} \sqrt{2 \log \left(\frac{2}{\delta} \right) + \log \det(M_{1,k+1})};$$

$$B_{h,k}^{MAT} \leftarrow \sqrt{|\mathcal{S}||\mathcal{A}|} + \frac{H-h+1}{2} \sqrt{2 \log \left(\frac{2}{\delta} \right) + \log \det(A_{k+1})};$$
 - 6: **if** $B_{h,k}^{VTR} \sqrt{X_{h,k}^\top M_{1,k+1}^{-1} X_{h,k}} \leq B_{h,k}^{MAT} \sqrt{\phi^\top(s, a) A_n^{-1} \phi(s, a)}$ **then**
 - 7: Perform one step of value iteration using the VTR model as follows: $Q_{h,k}(s, a) = r(s, a) + X_{h,k}^\top \theta_k + \sqrt{\beta_{h,k}} \sqrt{X_{h,k}^\top M_{1,k+1}^{-1} X_{h,k}}$
 - 8: **else**
 - 9: Update $Q_{h,k}(s, a)$ according to Equation 8 (Yang & Wang, 2019a) using the UC-MatrixRL model A_k . Note that in (Yang & Wang, 2019a) they use n to denote the current episode, in our paper we use k to denote the current episode.
 - 10: **end if**
 - 11: $V_{h,k}(s) = \max_a Q_{h,k}(s, a)$
 - 12: **end for**
 - 13: **end for**
-

We are now using multiple models instead of a single model, we must adjust our confidence sets accordingly. By using a union bound we replace δ with $\delta/2$ for our confidence parameter. This updated confidence parameter changes the term inside the logarithm. We now have $\log(2/\delta)$ where as before we had $\log(1/\delta)$.

G.2. Numerical Results

We will include the cumulative regret and the weighted L1 norm of UCRL-MIX on the RiverSwim environment as in Section 6. We also include a bar graph of the relative frequency with which the algorithm used the VTR-model for planning and the canonical model for planning.

Model-Based Reinforcement Learning with Value-Targeted Regression

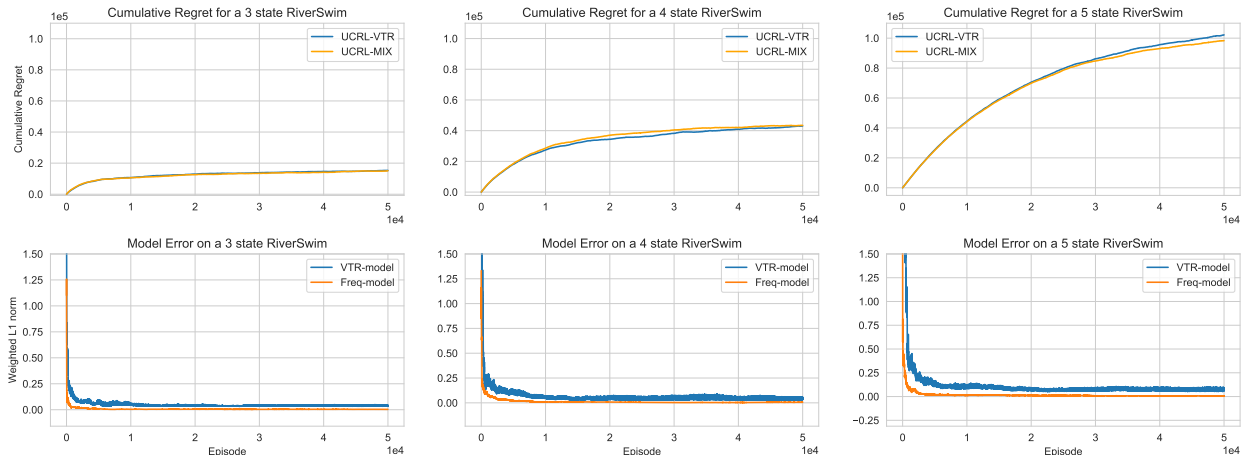


Figure 5. In the plots for the model error we include model error for both the VTR-model and the canonical model. Even though only one is used during planning both are updated at the end of each episode.

If we compare the results of Figure 5 with the results of Figure 2 from Section 6.2 we see that the cumulative regret of UCRL-MIX is almost identical to the cumulative regret of UCRL-VTR. The model errors of both the VTR and the canonical models are almost identical to the model errors of UCRL-VTR and UC-MatrixRL respectively.

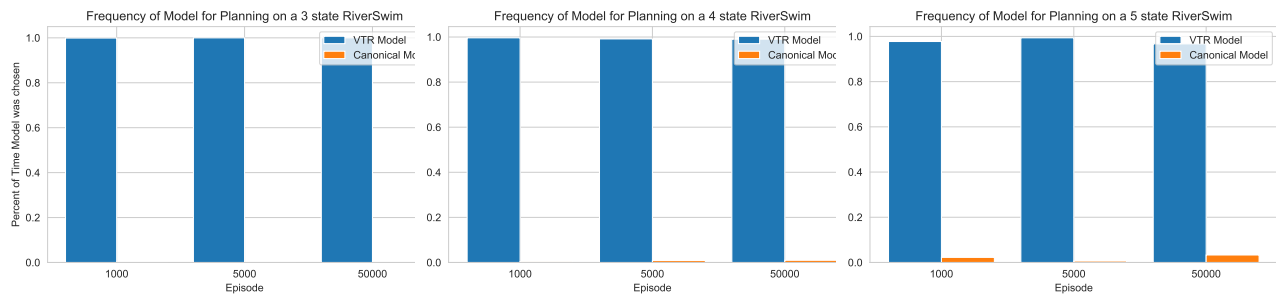


Figure 6. UCRL-MIX rarely, if ever, chooses the canonical model for planning on the RiverSwim environments.

From Figure 6, we see that on the RiverSwim environment, UCRL-MIX almost always uses the VTR-model for planning. We calculate this frequency by counting the number of times Step 7 of Alg 7 was observed up until episode k and by counting the number of times Step 9 of Alg 7 was observed up until episode k . We then divide these counts by the sum of the counts to get a percentage. We believe the reason the algorithm overwhelmingly chose the VTR-model was due to the fact that the confidence intervals for the VTR-model shrink much faster than the confidence intervals for the canonical model. The canonical model is forced to explore much longer than the VTR-model as its objective is to learn a globally optimal model rather than a model that yields high reward. Thus, the canonical model is forced to explore all state-action-next state tuples, even ones that do not yield high reward, in order to meet its objective of learning a globally optimal model while the VTR-model is only forced to explore state-action-next state tuples that fall in-line with its objective of accumulating high reward. The set of all state-action-next state tuples is much larger than the set of state-action-next state tuples that yield high reward which means the confidence intervals for the canonical model shrink slower than the confidence sets of the VTR-model on the RiverSwim environment.