

OPTIMISTIC POLICY OPTIMIZATION IS PROVABLY EFFICIENT IN NON-STATIONARY MDPs

Anonymous authors

Paper under double-blind review

ABSTRACT

We study episodic reinforcement learning (RL) in non-stationary linear kernel Markov decision processes (MDPs). In this setting, both the reward function and the transition kernel are linear with respect to the given feature maps and are allowed to vary over time, as long as their respective parameter variations do not exceed certain variation budgets. We propose the periodically restarted optimistic policy optimization algorithm (PROPO), which is an optimistic policy optimization algorithm with linear function approximation. PROPO features two mechanisms: sliding-window-based policy evaluation and periodic-restart-based policy improvement, which are tailored for policy optimization in a non-stationary environment. In addition, only utilizing the technique of sliding window, we propose a value-iteration algorithm. We establish dynamic upper bounds for the proposed methods and a matching minimax lower bound which shows the (near-) optimality of the proposed methods. To our best knowledge, PROPO is the first provably efficient policy optimization algorithm that handles non-stationarity.

1 INTRODUCTION

Reinforcement Learning (RL) (Sutton & Barto, 2018), coupled with powerful function approximators such as deep neural network, has demonstrated great potential in solving complicated sequential decision-making tasks such as games (Silver et al., 2016; 2017; Vinyals et al., 2019) and robotic control (Kober et al., 2013; Gu et al., 2017; Akkaya et al., 2019; Andrychowicz et al., 2020). Most of these empirical successes are driven by deep policy optimization methods such as trust region policy optimization (TRPO) (Schulman et al., 2015) and proximal policy optimization (PPO) (Schulman et al., 2017), whose performance has been extensively studied recently (Agarwal et al., 2019; Liu et al., 2019; Shani et al., 2020; Mei et al., 2020; Cen et al., 2020).

While classical RL assumes that an agent interacts with a time-invariant (stationary) environment, when deploying RL to real-world applications, both the reward function and Markov transition kernel can be time-varying. For example, in autonomous driving (Sallab et al., 2017), the vehicle needs to handle varying conditions of weather and traffic. When the environment changes with time, the agent must quickly adapt its policy to maximize the expected total rewards in the new environment. Meanwhile, another example of such a non-stationary scenario is when the environment is subject to adversarial manipulations, which is the case of adversarial attacks (Pinto et al., 2017; Huang et al., 2017; Pattanaik et al., 2017). In this situation, it is desired that the RL agent is robust against the malicious adversary.

Although there is a huge body of literature on developing provably efficient RL methods, most the existing works focus on the classical stationary setting, with a few exceptions include Jaksch et al. (2010); Gajane et al. (2018); Cheung et al. (2019a;c; 2020); Fei et al. (2020); Mao et al. (2020); Ortner et al. (2020); Domingues et al. (2020); Zhou et al. (2020b); Touati & Vincent (2020). However, these works all focus on value-based methods which only output greedy policies, and mostly focus on the tabular case where the state space is finite. Thus, the following problem remains open:

How can we design a provably efficient policy optimization algorithm for non-stationary environment in the context of function approximation?

There are four intertwined challenges associated with this problem: (i) bandit feedbacks from non-stationary reward and transition kernel, (ii) exploration-exploitation tradeoff that is inherent to online RL, (iii) incorporating function approximation in the algorithm, and (iv) characterizing the convergence and optimality of policy optimization. Existing works merely address a subset of these four challenges and it remains open how to tackle all of them simultaneously. For example, a line of research develops optimism-based value iteration algorithms that successfully handle (ii) and (iii), e.g., (Jiang et al., 2017; Jin et al., 2019b; Wang et al., 2019b; Zanette et al., 2020; Wang et al., 2020; Ayoub et al., 2020; Zhou et al., 2020a). Besides, Cai et al. (2019); Agarwal et al. (2020); Efroni et al. (2020) address challenges (ii)–(iv) but fail to consider (i), and Zhou et al. (2020b); Touati & Vincent (2020) tackle (i)–(iii) but leave (iv) open. More importantly, these four challenges are coupled together, which requires sophisticated algorithm design. In particular, due to challenges (i) and (iii), we need to track the non-stationary reward function and transition kernel by function estimation based on the bandit feedbacks. The estimated model is also time-varying and thus the corresponding policy optimization problem (challenge (iv)) has a non-stationary objective function. Moreover, to obtain sample efficiency, we need to strike a balance between exploration and exploitation in the policy update steps (challenge (i)).

In this work, we propose a periodically restarted optimistic policy optimization algorithm (PROPO) which successfully tackle the four challenges above. Specifically, we focus on the model of episodic linear kernel MDP (Ayoub et al., 2020; Zhou et al., 2020a) where both the reward and transition functions are parameterized by linear functions. Besides, we focus on the non-stationary setting and adopt the dynamic regret as the performance metric. Moreover, PROPO performs a policy evaluation step and a policy improvement step in each iteration. To handle challenges (i)–(iii), we propose a novel optimistic policy evaluation method that incorporates the technique of sliding window to handle non-stationarity. Specifically, based on the non-stationary bandit feedbacks, we propose to estimate the time-varying model via a sliding-window-based least-squares regression problem, where we only keep a subset of recent samples in regression. Based on the model estimator, we construct an optimistic value function by implementing model-based policy evaluation and adding an exploration bonus. Then, using such an optimistic value function as the update direction, in the policy improvement step, we propose to obtain a new policy by solving a Kullback-Leibler (KL) divergence regularized problem, which can be viewed as a mirror descent step. Moreover, as the underlying optimal policy is time-varying (challenge (iv)), we additionally restart the policy periodically by setting it to uniform policy every τ episodes. The two novel mechanisms, sliding window and periodic restart, respectively enable us to track the non-stationary MDP based on bandit feedbacks and handle the time-varying policy optimization problem.

Finally, to further exhibit effect of these two mechanisms, we propose an optimism-based value iteration algorithm, dubbed as SW-LSVI-UCB, which only utilize the sliding window and does not restart the policy as challenge (iv) disappears.

Our Contributions Our contribution is four-fold. First, we propose PROPO, a policy optimization algorithm designed for non-stationary linear kernel MDPs. This algorithm features two novel mechanisms, namely sliding window and periodic restart, and also incorporates linear function approximation and a bonus function to incentivize exploration. Second, we prove that PROPO achieves a sublinear dynamic regret, where d is the feature dimension, Δ is the total variation budget, H is the episode horizon, and T is the total number of steps. Third, to separately demonstrate the effect of sliding window, we propose a value-iteration algorithm, SW-LSVI-UCB, which adopts sliding-window-based regression to handle non-stationarity. Such an algorithm is shown to achieve a $\tilde{O}(d^{5/6}\Delta^{1/3}HT^{2/3})$ dynamic regret. Finally, we establish a $\Omega(d^{5/6}\Delta^{1/3}H^{2/3}T^{2/3})$ lower bound on the dynamic regret, which shows the (near-)optimality of the proposed algorithms. To our best knowledge, PROPO is the first provably efficient policy optimization algorithm under the non-stationary environment.

Related Work Our work adds to the vast body of existing literature on non-stationary MDPs. A line of work studies non-stationary RL in the tabular setting. See Jaksch et al. (2010); Gajane et al. (2018); Cheung et al. (2019a;c; 2020); Fei et al. (2020); Mao et al. (2020); Ortner et al. (2020) and the references therein for details. Recently, Domingues et al. (2020) consider the non-stationary RL in continuous environments and proposes a kernel-based algorithm. More related works are Zhou et al. (2020b); Touati & Vincent (2020), which study non-stationary linear MDPs, but their setting is not directly comparable with ours since linear MDPs cannot imply linear kernel MDPs. More-

over, Zhou et al. (2020b); Touati & Vincent (2020) do not incorporate policy optimization methods, which are more difficult because we need to handle the variation of the optimal policies of adjacent episodes and value-based methods only need to handle the non-stationarity drift of reward functions and transition kernels. Fei et al. (2020) also makes an attempt to investigate policy optimization algorithm for non-stationary environments. However, this work requires full-information feedback and only focuses on the tabular MDPs with time-varying reward functions and time-invariant transition kernels.

As a special case of MDP problems with unit horizon, bandit problems have been the subject of intense recent interest. See Besbes et al. (2014; 2019); Russac et al. (2019); Cheung et al. (2019a); Chen et al. (2019) and the references therein for details.

Another line closely related to our work is policy optimization. As proved in Yang et al. (2019); Agarwal et al. (2019); Liu et al. (2019); Wang et al. (2019a), policy optimization enjoys computational efficiency. Recently Cai et al. (2019); Efroni et al. (2020); Agarwal et al. (2020) proposed optimistic policy optimization methods which simultaneously attain computational efficiency and sample efficiency. Our work is also related to the value-based methods, especially LSVI (Bradtke & Barto, 1996; Jiang et al., 2017; Jin et al., 2019b; Wang et al., 2019b; Zanette et al., 2020; Wang et al., 2020; Ayoub et al., 2020; Zhou et al., 2020a).

Broadly speaking, our work is also related to a line of research on adversarial MDPs (Even-Dar et al., 2009; Neu et al., 2010; 2012; Zimin & Neu, 2013; Rosenberg & Mansour, 2019; Jin et al., 2019a).

Notation See §A for details.

2 PRELIMINARIES

2.1 NON-STATIONARY MDPs

An episodic non-stationary MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, H, P, r)$, where \mathcal{S} is a state space, \mathcal{A} is an action space, H is the length of each episode, $P = \{P_h^k\}_{(k,h) \in [K] \times [H]}$, $r = \{r_h^k\}_{(k,h) \in [K] \times [H]}$, where $P_h^k : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the probability transition kernel at the h -th step of the k -th episode, and $r_h^k : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the reward function at the h -th step of the k -th episode. We consider an agent which iteratively interacts with a non-stationary MDP in a sequence of K episodes. At the beginning of the k -th episode, the initial state s_1^k is adversarially given to the agent, and the agent determines a policy $\pi^k = \{\pi_h^k\}_{h=1}^H$. Then, at each step $h \in [H]$, the agent observes the state s_h^k , takes an action following the policy $a_h^k \sim \pi_h^k(\cdot | s_h^k)$ and receives a reward $r_h^k(s_h^k, a_h^k)$. Meanwhile, the MDP evolves into next state $s_{h+1}^k \sim P_h^k(\cdot | s_h^k, a_h^k)$. The k -th episode ends at state s_{H+1}^k , when this happens, no control action is taken and reward is equal to zero. We define the state and state-action value functions of policy $\pi = \{\pi_h\}_{h=1}^H$ recursively via the following Bellman equation:

$$Q_h^{\pi,k}(s, a) = r_h^k(s, a) + (\mathbb{P}_h^k V_{h+1}^{\pi,k})(s, a), \quad V_h^{\pi,k}(s) = \langle Q_h^{\pi,k}(s, \cdot), \pi_h(\cdot | s) \rangle_{\mathcal{A}}, \quad V_{H+1}^{\pi,k} = 0, \quad (2.1)$$

where \mathbb{P}_h^k is the operator form of the transition kernel $P_h^k(\cdot | \cdot, \cdot)$, which is defined as

$$(\mathbb{P}_h^k f)(s, a) = \mathbb{E}[f(s') | s' \sim P_h^k(s' | s, a)] \quad (2.2)$$

for any function $f : \mathcal{S} \rightarrow \mathbb{R}$. Here $\langle \cdot, \cdot \rangle_{\mathcal{A}}$ denotes the inner product over \mathcal{A} .

In the literature of optimization and reinforcement learning, the performance of the agent is measured by its dynamic regret, which measures the difference between the agent's policy and the benchmark policy $\pi^* = \{\pi^{*,k}\}_{k=1}^K$. Specifically, the dynamic regret is defined as

$$\text{D-Regret}(T, \pi^*) = \sum_{k=1}^K (V_1^{\pi^{*,k},k}(s_1^k) - V_1^{\pi^k,k}(s_1^k)), \quad (2.3)$$

where $T = HK$ is the number of steps taken by agent and $\pi^{*,k}$ is the benchmark policy of episode k . It is worth mentioning that when the benchmark policy is the optimal policy of each individual episode, that is, $\pi^{*,k} = \arg\max_{\pi} V_1^{\pi,k}(s_1^k)$, the dynamic regret reaches the maximum, and this special case is widely considered in previous works (Cheung et al., 2020; Mao et al., 2020; Domingues et al., 2020). Throughout the paper, when π^* is clear from the context, we may omit π^* from $\text{D-Regret}(T, \pi^*)$.

2.2 MODEL ASSUMPTIONS

We focus on the linear setting of Markov decision process, where the reward functions and transition kernels are assumed to be linear. We formally make the following assumption.

Assumption 2.1 (Non-stationary Linear Kernel MDP). MDP $(\mathcal{S}, \mathcal{A}, H, P, r)$ is a linear kernel MDP with known feature maps $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and $\psi : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^d$, if for any $(k, h) \in [K] \times [H]$, there exist unknown vectors $\theta_h^k \in \mathbb{R}^d$ and $\xi_h^k \in \mathbb{R}^d$, such that

$$r_h^k(s, a) = \phi(s, a)^\top \theta_h^k, \quad P_h^k(s' | s, a) = \psi(s, a, s')^\top \xi_h^k$$

for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$. Without loss of generality, we assume that

$$\|\phi(s, a)\|_2 \leq 1, \quad \|\theta_h^k\|_2 \leq \sqrt{d}, \quad \|\xi_h^k\|_2 \leq \sqrt{d}$$

for any $(k, h) \in [K] \times [H]$. Moreover, we assume that

$$\int_{\mathcal{S}} \|\psi(s, a, s')\|_2 ds' \leq \sqrt{d}$$

for any $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Our assumption consists of two parts. One is about reward functions, which follows the setting of linear bandits (Abbasi-Yadkori et al., 2011; Agrawal & Goyal, 2013; Besbes et al., 2014; 2015; Cheung et al., 2019a;b). The other part is about transition kernels. As shown in Cai et al. (2019); Ayoub et al. (2020); Zhou et al. (2020a), linear kernel MDPs as defined above cover several other MDPs studied in previous works, as special cases. For example, tabular MDPs with canonical basis (Cai et al., 2019; Ayoub et al., 2020; Zhou et al., 2020a), feature embedding of transition models (Yang & Wang, 2019a) and linear combination of base models (Modi et al., 2020) are special cases. However, it is worth mentioning that Jin et al. (2019b); Yang & Wang (2019b) studied another “linear MDPs”, which assumes the transition kernels can be represented as $\mathbb{P}_h(s' | s, a) = \psi'(s, a)^\top \mu_h(s')$ for any $h \in [H]$ and $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$. Here $\psi'(\cdot, \cdot)$ is a known feature map and $\mu(\cdot)$ is an unknown measure. It is worth noting that linear MDPs studied in our paper and linear MDPs (Jin et al., 2019b; Yang & Wang, 2019b) are two different classes of MDPs since their feature maps $\psi(\cdot, \cdot, \cdot)$ and $\psi'(\cdot, \cdot)$ are different and neither class of MDPs includes the other.

To facilitate the following analysis, we denote by $P_h^{k, \pi}$ the Markov kernel of policy at the h -step of the k -th episode, that is, for $s \in \mathcal{S}$, $P_h^{k, \pi}(\cdot | s) = \sum_{a \in \mathcal{A}} P_h^k(\cdot | s, a) \cdot \pi_h(a | s)$. Also, we define

$$\begin{aligned} \|\pi_h - \pi'_h\|_{\infty, 1} &= \max_{s \in \mathcal{S}} \|\pi_h(\cdot | s) - \pi'_h(\cdot | s)\|_1, \\ \|P_h^{k, \pi} - P_h^{k, \pi'}\|_{\infty, 1} &= \max_{s \in \mathcal{S}} \|P_h^{k, \pi}(\cdot | s) - P_h^{k, \pi'}(\cdot | s)\|_1. \end{aligned}$$

Next, we introduce several measures of change in MDPs. First, we denote by P_T the total variation in the benchmark policies of adjacent episodes:

$$P_T = \sum_{k=1}^K \sum_{h=1}^H \|\pi_h^{*, k} - \pi_h^{*, k-1}\|_{\infty, 1}, \quad (2.4)$$

where we choose $\pi_h^{*, 0} = \pi_h^{*, 1}$ for any $h \in [H]$.

Next, we assume the drifting environment (Besbes et al., 2014; 2015; Cheung et al., 2019a; Russac et al., 2019), that is, θ_h^k and ξ_h^k can change over different indexes (k, h) , with the constraint that the sum of the Euclidean distances between consecutive θ_h^k and ξ_h^k are bounded by variation budgets B_T and B_P , that is,

$$\sum_{h=1}^H \sum_{k=1}^K \|\theta_h^{k-1} - \theta_h^k\|_2 \leq B_T, \quad \sum_{h=1}^H \sum_{k=1}^K \|\xi_h^{k-1} - \xi_h^k\|_2 \leq B_P, \quad \Delta = B_T + B_P, \quad (2.5)$$

where H is the length of each episode, K is the total number of episodes, and $T = HK$ is the total number of steps taken by the agent. Here Δ is the total variation budget, which quantifies the non-stationarity of a linear kernel MDP.

3 MINIMAX LOWER BOUND

In this section, we provide the information-theoretical lower bound result. The following theorem shows a minimax lower bound of dynamic regret for any algorithm to learn non-stationary linear kernel MDPs.

Theorem 3.1 (Minimax lower bound). Fix $\Delta > 0$, $H > 0$, $d \geq 2$, and $T = \Omega(d^{5/2} \Delta H^{1/2})$. Then, there exists a non-stationary linear kernel MDP with a d -dimensional feature map and maximum total variation budget Δ , such that,

$$\min_{\mathbb{A}} \max_{\pi^*} \text{D-Regret}(T, \pi^*) \geq \Omega(d^{5/6} \Delta^{1/3} H^{2/3} T^{2/3}),$$

where \mathbb{A} denotes the learning algorithm.

Proof sketch. As mentioned above, we only need to establish the lower bound of the dynamic regret when the benchmark policy is the optimal policy of each individual episode. The proof of lower bound relies on the construction of a hard-to-learn non-stationary linear kernel MDP instance. To handle the non-stationarity, we need to divide the total T steps into L segments, where each segment contains $T_0 = \lfloor \frac{T}{L} \rfloor$ steps and has $K_0 = \lfloor \frac{K}{L} \rfloor$ episodes. Within each segment, the construction of MDP is similar to the hard-to-learn instance constructed in stationary RL problems (Jaksch et al., 2010; Lattimore & Hutter, 2012; Osband & Van Roy, 2016). Then, we can derive a lower bound of $\Omega(dH\sqrt{T_0})$ for the stationary RL problem. Meanwhile, the transition kernel of this hard-to-learn MDP changes abruptly between two consecutive segments, which forces the agent to learn a new stationary MDP in each segment. Finally, by optimization L subject to the total budget constraint, we obtain the lower bound of $\Omega(d^{5/6} \Delta^{1/3} H^{2/3} T^{2/3})$. See Appendix C for details. \square

4 ALGORITHM AND THEORY

4.1 PROPO

Now we present Periodically Restarted Optimistic Policy Optimization (PROPO) in Algorithm 1, which includes a policy improvement step and a policy evaluation step.

Policy Improvement Step. At k -th episode, Model-Based OPPO updates $\pi^k = \{\pi_h^k\}_{h=1}^H$ according to $\pi^{k-1} = \{\pi_h^{k-1}\}_{h=1}^H$. Motivated by the policy improvement step in NPG (Kakade, 2002), TRPO (Schulman et al., 2015), and PPO (Schulman et al., 2017), we consider the following policy improvement step

$$\pi^k = \operatorname{argmax}_{\pi} L_{k-1}(\pi), \quad (4.1)$$

where $L_{k-1}(\pi)$ is defined as

$$\begin{aligned} L_{k-1}(\pi) = \mathbb{E}_{\pi^{k-1}} \left[\sum_{h=1}^H \langle Q_h^{k-1}(s_h, \cdot), \pi_h(\cdot | s_h) - \pi_h^{k-1}(\cdot | s_h) \rangle \right] \\ - \alpha^{-1} \cdot \mathbb{E}_{\pi^{k-1}} \left[\sum_{h=1}^H \text{KL}(\pi_h(\cdot | s_h) \parallel \pi_h^{k-1}(\cdot | s_h)) \right], \end{aligned} \quad (4.2)$$

where $\alpha > 0$ is a stepsize and Q_h^{k-1} which is obtained in Line 10 of Algorithm 2 is the estimator of $Q_h^{\pi^{k-1}, k-1}$. Here the expectation $\mathbb{E}_{\pi^{k-1}}$ is taken over the random state-action pairs $\{(s_h, a_h)\}_{h=1}^H$, where the initial state $s_1 = s_1^k$, the distribution of action a_h follows $\pi(\cdot | s_h)$, and the distribution of the next state s_{h+1} follows the transition dynamics $P_h^k(\cdot | s_h, a_h)$. Such a policy improvement step can also be regarded as one iteration of infinite-dimensional mirror descent (Nemirovsky & Yudin, 1983; Liu et al., 2019; Wang et al., 2019a).

By the optimality condition, policy update in (4.1) admits a closed-form solution

$$\pi_h^k(\cdot | s) \propto \pi_h^{k-1}(\cdot | s) \cdot \exp\{\alpha \cdot Q_h^{k-1}(s, \cdot)\} \quad (4.3)$$

for any $s \in \mathcal{S}$ and $(k, h) \in [K] \times [H]$.

Policy Evaluation Step. At the end of the k -th episode, Model-Based OPPO evaluates the policy π^k based on the $(k-1)$ historical trajectories. Then, we show the details of estimating the reward functions and transition kernels, respectively.

(i) Estimating Reward. To estimate the reward functions, we use the sliding window regularized least squares estimator (SW-RLSE) (Garivier & Moulines, 2011; Cheung et al., 2019a;b), which is a key tool in estimating the unknown parameters online. At h -th step of k -th episode, we aim to estimate the unknown parameter θ_h^k based on the historical observation $\{(s_h^\tau, a_h^\tau), r_h^\tau(s_h^\tau, a_h^\tau)\}_{\tau=1}^{k-1}$. The design of SW-RLSE is based on the “forgetting principle” (Garivier & Moulines, 2011), that is, under non-stationarity, the historical observations far in the past are obsolete, and they do not contain relevant information for the evaluation of the current policy. Therefore, we could estimate θ_h^k using only $\{(s_h^\tau, a_h^\tau), r_h^\tau(s_h^\tau, a_h^\tau)\}_{\tau=1 \vee (k-w)}^{k-1}$, the observations during the sliding window $1 \vee (k-w)$ to $k-1$,

$$\hat{\theta}_h^k = \underset{\theta}{\operatorname{argmin}} \left(\sum_{\tau=1 \vee (k-w)}^{k-1} (r_h^\tau(s_h^\tau, a_h^\tau) - \phi(s_h^k, a_h^k)^\top \theta)^2 + \lambda \cdot \|\theta\|_2^2 \right), \quad (4.4)$$

where λ is the regularization parameter and w is the length of a sliding window. By solving (4.4), we obtain the estimator of θ_h^k :

$$\begin{aligned} \hat{\theta}_h^k &= (\Lambda_h^k)^{-1} \sum_{\tau=1 \vee (k-w)}^{k-1} \phi(s_h^\tau, a_h^\tau) r_h^\tau(s_h^\tau, a_h^\tau), \\ \text{where } \Lambda_h^k &= \sum_{\tau=1 \vee (k-w)}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda I_d. \end{aligned} \quad (4.5)$$

(ii) Estimating Transition. Similar to the estimation of reward functions, for any $(k, h) \in [K] \times [H]$, we define the sliding window empirical mean-squared Bellman error (SW-MSBE) as

$$M_h^k(\xi) = \sum_{\tau=1 \vee (k-w)}^{k-1} (V_{h+1}^\tau(s_{h+1}^\tau) - \eta_h^\tau(s_h^\tau, a_h^\tau)^\top \xi)^2,$$

where we denote $\eta_h^\tau(\cdot, \cdot)$ as

$$\eta_h^\tau(\cdot, \cdot) = \int_{\mathcal{S}} \psi(\cdot, \cdot, s') \cdot V_{h+1}^\tau(s') ds'. \quad (4.6)$$

By Assumption 2.1, we have

$$\|\eta_h^k(\cdot, \cdot)\|_2 \leq H\sqrt{d}$$

for any $(k, h) \in [K] \times [H]$. Then we estimate ξ_h^k by solving the following problem:

$$\hat{\xi}_h^k = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} (M_h^k(w) + \lambda' \cdot \|w\|_2^2), \quad (4.7)$$

where λ' is the regularization parameter. By solving (4.7), we obtain

$$\begin{aligned} \hat{\xi}_h^k &= (A_h^k)^{-1} \left(\sum_{\tau=1 \vee (k-w)}^{k-1} \eta_h^\tau(s_h^\tau, a_h^\tau) \cdot V_{h+1}^\tau(s_{h+1}^\tau) \right) \\ \text{where } A_h^k &= \sum_{\tau=1 \vee (k-w)}^{k-1} \eta_h^\tau(s_h^\tau, a_h^\tau) \eta_h^\tau(s_h^\tau, a_h^\tau)^\top + \lambda' I_d. \end{aligned} \quad (4.8)$$

The policy evaluation step is iteratively updating the estimated Q-function $Q^k = \{Q_h^k\}_{h=1}^H$ by

$$\begin{aligned} Q_h^k(\cdot, \cdot) &= \min\{\phi(\cdot, \cdot)^\top \hat{\theta}_h^k + \eta_h^k(\cdot, \cdot)^\top \hat{\xi}_h^k + B_h^k(\cdot, \cdot) + \Gamma_h^k(\cdot, \cdot), H - h + 1\}^+, \\ V_h^k(s) &= \langle Q_h^k(s, \cdot), \pi_h^k(\cdot|s) \rangle_{\mathcal{A}} \end{aligned} \quad (4.9)$$

in the order of $h = H, H - 1, \dots, 1$. Here bonus functions $B_h^k(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$ and $\Gamma_h^k(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$ are used to quantify the uncertainty in estimating reward r_h^k and quantity $\mathbb{P}_h^k V_{h+1}^k$ respectively, defined as

$$B_h^k(\cdot, \cdot) = \beta(\phi(\cdot, \cdot)^\top (\Lambda_h^k)^{-1} \phi(\cdot, \cdot))^{1/2}, \quad \Gamma_h^k(\cdot, \cdot) = \beta'(\eta_h^k(\cdot, \cdot)^\top (A_h^k)^{-1} \eta_h^k(\cdot, \cdot))^{1/2}, \quad (4.10)$$

where $\beta > 0$ and $\beta' > 0$ are parameters depend on d, H and K , which are specified in Theorem 4.2.

To handle the non-stationary drift incurred by the different optimal policies in different episodes, Algorithm 1 also includes a periodic restart mechanism, which resets the policy estimates every τ episodes. We call the τ episodes between every two resets a segment. In each segment, each episode is approximately the same as the first episode, which means that we can regard it as a stationary MDP. Then we can use the method of solving the stationary MDP to analyze each segment with a small error, and finally combine each segment and choose the value of τ to get the desired result. Such a restart mechanism is widely used in RL (Auer et al., 2009; Ortner et al., 2020), bandits (Besbes et al., 2014; Zhao et al., 2020), and non-stationary optimization (Besbes et al., 2015; Jadbabaie et al., 2015).

The pseudocode of the PROPO algorithm is given in Algorithm 1.

Algorithm 1 Periodically Restarted Optimistic Policy Optimization (PROPO)

Require: Reset cycle length τ , sliding window length w , stepsize α , regularization factors λ and λ' , and bonus multipliers β and β' .

- 1: Initialize $\{\pi_h^0(\cdot | \cdot)\}_{h=1}^H$ as uniform distribution policies, $\{Q_h^0(\cdot, \cdot)\}_{h=1}^H$ as zero functions.
- 2: **for** $k = 1, 2, \dots, K$ **do**
- 3: Receive the initial state s_1^k .
- 4: **if** $k \bmod \tau = 1$ **then**
- 5: Set $\{Q_h^{k-1}\}_{h \in [H]}$ as zero functions and $\{\pi_h^{k-1}\}_{h \in [H]}$ as uniform distribution on \mathcal{A} .
- 6: **end if**
- 7: **for** $h = 1, 2, \dots, H$ **do**
- 8: $\pi_h^k(\cdot | \cdot) \propto \pi_h^{k-1}(\cdot | \cdot) \cdot \exp\{\alpha \cdot Q_h^{k-1}(\cdot, \cdot)\}$.
- 9: Take action $a_h^k \sim \pi_h^k(\cdot | s_h^k)$.
- 10: Observe the reward $r_h^k(s_h^k, a_h^k)$ and receive the next state s_{h+1}^k .
- 11: **end for**
- 12: Compute Q_h^k by SWOPE($k, \{\pi_h^k\}, \lambda, \lambda', \beta, \beta'$) (Algorithm 2).
- 13: **end for**

Algorithm 2 Sliding Window Optimistic Policy Evaluation (SWOPE)

Require: Episode index k , policies $\{\pi_h\}$, regularization factors λ and λ' , and bonus multipliers β and β' .

- 1: Initialize V_{H+1}^k as a zero function.
- 2: **for** $h = H, H - 1, \dots, 0$ **do**
- 3: $\eta_h^k(\cdot, \cdot) = \int_{\mathcal{S}} \psi(\cdot, \cdot, s') \cdot V_{h+1}^k(s') ds'$.
- 4: $\Lambda_h^k = \sum_{\tau=1 \vee (k-w)}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda I_d$.
- 5: $\hat{\theta}_h^k = (\Lambda_h^k)^{-1} \sum_{\tau=1 \vee (k-w)}^{k-1} \phi(s_h^\tau, a_h^\tau) r_h^\tau(s_h^\tau, a_h^\tau)$.
- 6: $A_h^k = \sum_{\tau=1 \vee (k-w)}^{k-1} \eta_h^\tau(s_h^\tau, a_h^\tau) \eta_h^\tau(s_h^\tau, a_h^\tau)^\top + \lambda' I_d$.
- 7: $\hat{\xi}_h^k = (A_h^k)^{-1} (\sum_{\tau=1 \vee (k-w)}^{k-1} \eta_h^\tau(s_h^\tau, a_h^\tau) \cdot V_{h+1}^\tau(s_{h+1}^\tau))$.
- 8: $B_h^k(\cdot, \cdot) = \beta(\phi(\cdot, \cdot)^\top (\Lambda_h^k)^{-1} \phi(\cdot, \cdot))^{1/2}$.
- 9: $\Gamma_h^k(\cdot, \cdot) = \beta'(\eta_h^k(\cdot, \cdot)^\top (A_h^k)^{-1} \eta_h^k(\cdot, \cdot))^{1/2}$.
- 10: $Q_h^k(\cdot, \cdot) = \min\{\phi(\cdot, \cdot)^\top \hat{\theta}_h^k + \eta_h^k(\cdot, \cdot)^\top \hat{\xi}_h^k + B_h^k(\cdot, \cdot) + \Gamma_h^k(\cdot, \cdot), H - h + 1\}^+$.
- 11: $V_h^k(s) = \langle Q_h^k(s, \cdot), \pi_h^k(\cdot | s) \rangle_{\mathcal{A}}$.
- 12: **end for**

4.2 SW-LSVI-UCB

In this subsection, we present the details of Sliding Window Least-Square Value Iteration with UCB (SW-LSVI-UCB) in Algorithm 3 (cf. Appendix B).

Similar to Least-Square Value Iteration with UCB (LSVI-UCB) in Jin et al. (2019b), SW-LSVI-UCB is also an optimistic modification of Least-Square Value Iteration (LSVI) (Bradtke & Barto, 1996), where the optimism is realized by Upper-Confidence Bounds (UCB). Specifically, the optimism is achieved due to the bonus functions B_h^k and Γ_h^k , which quantify the uncertainty of reward functions and transition kernels, respectively. It is worth noting that in order to handle the non-stationarity, SW-LSVI-UCB also uses the sliding window method (Garivier & Moulines, 2011; Cheung et al., 2019a;b).

In detail, at k -th episode, SW-LSVI-UCB consists of two steps. In the first step, by solving the sliding window least-square problems (4.4) and (4.7), SW-LSVI-UCB updates the parameters Λ_h^k in (4.5), $\hat{\theta}_h^k$ in (4.5), A_h^k in (4.8), and $\hat{\xi}_h^k$ in (4.8), which are used to form the Q-function Q_h^k . In the second step, SW-LSVI-UCB obtains the greedy policy with respect to the Q-function Q_h^k gained in the first step. See Algorithm 3 in Appendix B for more details.

4.3 REGRET ANALYSIS

In this subsection, we analyze the dynamic regret incurred by Algorithms 1 and 3 and compare the theoretical regret upper bounds derived for these two algorithms.

To derive sharp dynamic regret bounds, we impose the following technical assumption.

Assumption 4.1. There exists an orthonormal basis $\Psi = (\Psi_1, \dots, \Psi_d)$ such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists a vector $z \in \mathbb{R}^d$ satisfying that $\phi(s, a) = \Psi z$. We also assume the existence of another orthonormal basis $\Psi' = (\Psi'_1, \dots, \Psi'_d)$ such that for any $(s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$ such that $\eta_h^k(s, a) = \Psi' z'$ for some $z' \in \mathbb{R}^d$.

It is not difficult to show that this assumption holds in the tabular setting. Similar assumption is also adopted by previous work in non-stationary optimization (Cheung et al., 2019a). We will provide more comments on this technical assumption after showing main results.

First, we establish an upper bound on the dynamic regret of PROPO. Recall that the dynamic regret is defined in (2.3) and d is the dimension of the feature maps ϕ and ψ . Also, $|\mathcal{A}|$ is the cardinality of \mathcal{A} . We also define $\rho = \lceil K/\tau \rceil$ to be the number of restarts that take place in Algorithm 1.

Theorem 4.2 (Upper bound for Algorithm 1). Suppose Assumptions 2.1 and 4.1 hold. Let $\tau = \Pi_{[1, K]}(\lfloor (\frac{T\sqrt{\log |\mathcal{A}|}}{H(P_T + \sqrt{d}\Delta)})^{2/3} \rfloor)$, $\alpha = \sqrt{\rho \log |\mathcal{A}| / (H^2 K)}$ in (4.2), $w = \Theta(d^{1/3} \Delta^{-2/3} T^{2/3})$ in (4.4), $\lambda = \lambda' = 1$ in (4.4) and (4.9), $\beta = \sqrt{d}$ in (4.10), and $\beta' = C' \sqrt{dH^2 \cdot \log(dT/\zeta)}$ in (4.10), where $C' > 1$ is an absolute constant and $\zeta \in (0, 1]$. We have

$$\begin{aligned} \text{D-Regret}(T) &\lesssim d^{5/6} \Delta^{1/3} H T^{2/3} \cdot \log(dT/\zeta) \\ &+ \begin{cases} \sqrt{H^3 T \log |\mathcal{A}|}, & \text{if } 0 \leq P_T + \sqrt{d}\Delta \leq \sqrt{\frac{\log |\mathcal{A}|}{K}}, \\ (H^2 T \sqrt{\log |\mathcal{A}|})^{2/3} (P_T + \sqrt{d}\Delta)^{1/3}, & \text{if } \sqrt{\frac{\log |\mathcal{A}|}{K}} \leq P_T + \sqrt{d}\Delta \lesssim K \sqrt{\log |\mathcal{A}|}, \\ H^2 (P_T + \sqrt{d}\Delta), & \text{if } P_T + \sqrt{d}\Delta \gtrsim K \sqrt{\log |\mathcal{A}|}, \end{cases} \end{aligned}$$

with probability at least $1 - \zeta$.

Proof. See Appendix D for a proof sketch and Appendix G for a detailed proof. \square

Then we discuss the regret bound throughout three regimes of $P_T + \sqrt{d}\Delta$:

- Small $P_T + \sqrt{d}\Delta$: when $0 \leq P_T + \sqrt{d}\Delta \leq \sqrt{\frac{\log |\mathcal{A}|}{K}}$, restart period $\tau = K$, which means that we do not need to periodically restart in this case. Assuming that $\log |\mathcal{A}| = \mathcal{O}(d^{5/3} \Delta^{2/3} H^{-1} T^{1/3})$, Algorithm 1 attains a $\tilde{\mathcal{O}}(d^{5/6} \Delta^{1/3} H T^{2/3})$ dynamic regret. Combined with the lower bound established in Theorem 3.1, our result matches the lower bound

in d , Δ and T up to logarithmic factors. Hence, we can conclude that Algorithm 1 is a near-optimal algorithm;

- Moderate $P_T + \sqrt{d}\Delta$: when $\sqrt{\frac{\log |\mathcal{A}|}{K}} \leq P_T + \sqrt{d}\Delta \lesssim K\sqrt{\log |\mathcal{A}|}$, restart period $\tau = (\frac{T\sqrt{\log |\mathcal{A}|}}{H(P_T + \sqrt{d}\Delta)})^{2/3} \in [2, K]$. Algorithm 2 incurs a $\tilde{\mathcal{O}}(T^{2/3})$ dynamic regret if $\Delta = \mathcal{O}(1)$ and $P_T = \mathcal{O}(1)$;
- Large $P_T + \sqrt{d}\Delta$: when $P_T + \sqrt{d}\Delta \gtrsim K\sqrt{\log |\mathcal{A}|}$, restart period $\tau = K$. Since the model is highly non-stationary, we only obtain a linear regret in T .

In the following theorem, we establish the upper bound of dynamic regret incurred by SW-LSVI-UCB (Algorithm 3).

Theorem 4.3 (Upper bound for Algorithm 3). Suppose Assumption 2.1 and 4.1 hold. Let $w = \Theta(d^{1/3}\Delta^{-2/3}T^{2/3})$ in (4.4), $\lambda = \lambda' = 1$ in (4.4) and (4.9), $\beta = \sqrt{d}$ in (4.10), and $\beta' = C'\sqrt{dH^2 \cdot \log(dT/\zeta)}$ in (4.10), where $C' > 1$ is an absolute constant and $\zeta \in (0, 1]$. We have

$$\text{D-Regret}(T) \lesssim d^{5/6}\Delta^{1/3}HT^{2/3} \cdot \log(dT/\zeta)$$

with probability at least $1 - \zeta$.

Proof. See Appendix H for a detailed proof. \square

Regarding Assumption 4.1. Due to some technical issue (Touati & Vincent, 2020; Zhao & Zhang, 2021), without this assumption and the knowledge of locally variation budget (Touati & Vincent, 2020), previous work can only obtain the bound $\tilde{\mathcal{O}}(T^{3/4})$ (Cheung et al., 2020; Zhao & Zhang, 2021; Zhao et al., 2020; Russac et al., 2019; Zhou et al., 2020b; Touati & Vincent, 2020). Thanks to Assumption 4.1, we derive sharper regret bounds at the order $\tilde{\mathcal{O}}(T^{2/3})$. We also remark that we can establish slightly worse regret bounds for Algorithms 1 and 3 without Assumption 4.1. See Appendix I for details.

Optimality of the Bounds. Notably, the term $\tilde{\mathcal{O}}(d^{5/6}\Delta^{1/3}HT^{2/3})$ appears in both the results in Theorems 4.2 and 4.3. Ignoring logarithmic factors, there is only a gap of $H^{1/3}$ between this upper bound and the lower bound $\Omega(d^{5/6}\Delta^{1/3}H^{2/3}T^{2/3})$ established in Theorem 3.1. We conjecture that this gap can be bridged by using the ‘‘Bernstein’’ type bonus functions Azar et al. (2017); Jin et al. (2018). Since our focus is on designing a provably efficient policy optimization algorithm for non-stationary linear kernel MDPs, we don’t use this technique for the clarity of our analysis.

Comparison. Compared with PROPO, SW-LSVI-UCB achieves a slightly better regret without the help of the periodic restart mechanism. Especially in the highly non-stationary case, that is $P_T + \sqrt{d}\Delta \gtrsim K\sqrt{\log |\mathcal{A}|}$, SW-LSVI-UCB achieves a $\tilde{\mathcal{O}}(T^{2/3})$ regret, where PROPO only attains a linear regret in T . However, PROPO achieves the same $\tilde{\mathcal{O}}(T^{2/3})$ regret as SW-LSVI-UCB when $P_T + \sqrt{d}\Delta \lesssim K\sqrt{\log |\mathcal{A}|}$, which suggests that PROPO is provably efficient for solving slightly or even moderately non-stationary MDPs. Therefore, it is important to investigate whether it is possible to bridge this gap between policy and value based methods, or alternatively to show that this gap is actually a true drawback of policy optimization methods in the non-stationary case.

5 CONCLUSION

In this work, we have proposed a probably efficient policy optimization algorithm, dubbed as PROPO, for non-stationary linear kernel MDPs. Such an algorithm incorporates a bonus function to incentivize exploration, and more importantly, adopts sliding-window-based regression in policy evaluation and periodic restart in policy update to handle the challenge of non-stationarity. Moreover, as a byproduct, we establish an optimistic value iteration algorithm, SW-LSVI-UCB, by combining UCB and sliding-window. We prove that PROPO and SW-LSVI-UCB both achieve sample efficiency by having sublinear dynamic regret. We also establish a dynamic regret lower bound which shows that PROPO and SW-LSVI-UCB are near-optimal. To our best knowledge, we propose the first provably efficient policy optimization method that successfully handles non-stationarity.

REFERENCES

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *arXiv preprint arXiv:1908.00261*, 2019.
- Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. PC-PG: Policy cover directed exploration for provable policy gradient learning. *arXiv preprint arXiv:2007.08459*, 2020.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pp. 127–135, 2013.
- Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. In *Advances in neural information processing systems*, pp. 89–96, 2009.
- Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin F Yang. Model-based reinforcement learning with value-targeted regression. *arXiv preprint arXiv:2006.01107*, 2020.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pp. 263–272. PMLR, 2017.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Stochastic multi-armed-bandit problem with non-stationary rewards. In *Advances in neural information processing systems*, pp. 199–207, 2014.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Non-stationary stochastic optimization. *Operations research*, 63(5):1227–1244, 2015.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Optimal exploration–exploitation in a multi-armed bandit problem with non-stationary rewards. *Stochastic Systems*, 9(4):319–337, 2019.
- Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*, 2019.
- Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *arXiv preprint arXiv:2007.06558*, 2020.
- Yifang Chen, Chung-Wei Lee, Haipeng Luo, and Chen-Yu Wei. A new algorithm for non-stationary contextual bandits: Efficient, optimal and parameter-free. In *Conference on Learning Theory*, pp. 696–726. PMLR, 2019.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Hedging the drift: Learning to optimize under non-stationarity. *arXiv preprint arXiv:1903.01461*, 2019a.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Learning to optimize under non-stationarity. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1079–1087, 2019b.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Non-stationary reinforcement learning: The blessing of (more) optimism. *Available at SSRN 3397818*, 2019c.

- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Reinforcement learning for non-stationary markov decision processes: The blessing of (more) optimism. *arXiv preprint arXiv:2006.14389*, 2020.
- Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. 2008.
- Omar Darwiche Domingues, Pierre Ménard, Matteo Pirota, Emilie Kaufmann, and Michal Valko. A kernel-based approach to non-stationary reinforcement learning in metric spaces. *arXiv preprint arXiv:2007.05078*, 2020.
- Yonathan Efroni, Lior Shani, Aviv Rosenberg, and Shie Mannor. Optimistic policy optimization with bandit feedback. *arXiv preprint arXiv:2002.08243*, 2020.
- Eyal Even-Dar, Sham M Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Yingjie Fei, Zhuoran Yang, Zhaoran Wang, and Qiaomin Xie. Dynamic regret of policy optimization in non-stationary environments. *arXiv preprint arXiv:2007.00148*, 2020.
- Pratik Gajane, Ronald Ortner, and Peter Auer. A sliding-window algorithm for markov decision processes with arbitrarily changing rewards and transitions. *arXiv preprint arXiv:1805.10066*, 2018.
- Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for switching bandit problems. In *International Conference on Algorithmic Learning Theory*, pp. 174–188. Springer, 2011.
- Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE international conference on robotics and automation (ICRA)*, pp. 3389–3396. IEEE, 2017.
- Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *arXiv preprint arXiv:1702.02284*, 2017.
- Ali Jadbabaie, Alexander Rakhlin, Shahin Shahrampour, and Karthik Sridharan. Online optimization: Competing with dynamic comparators. In *Artificial Intelligence and Statistics*, pp. 398–406, 2015.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pp. 1704–1713. PMLR, 2017.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pp. 4863–4873, 2018.
- Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial mdps with bandit feedback and unknown transition. *arXiv preprint arXiv:1912.01192*, 2019a.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019b.
- Sham M Kakade. A natural policy gradient. In *Advances in neural information processing systems*, pp. 1531–1538, 2002.
- Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Tor Lattimore and Marcus Hutter. Pac bounds for discounted mdps. In *International Conference on Algorithmic Learning Theory*, pp. 320–334. Springer, 2012.

- Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural trust region/proximal policy optimization attains globally optimal policy. In *Advances in Neural Information Processing Systems*, pp. 10565–10576, 2019.
- Weichao Mao, Kaiqing Zhang, Ruihao Zhu, David Simchi-Levi, and Tamer Başar. Near-optimal regret bounds for model-free rl in non-stationary episodic mdps, 2020.
- Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pp. 6820–6829. PMLR, 2020.
- Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*, pp. 2010–2020, 2020.
- Arkadii Semenovich Nemirovsky and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.
- Gergely Neu, András György, and Csaba Szepesvári. The online loop-free stochastic shortest-path problem. In *COLT*, volume 2010, pp. 231–243. Citeseer, 2010.
- Gergely Neu, Andras Gyorgy, and Csaba Szepesvári. The adversarial stochastic shortest path problem with unknown transition probabilities. In *Artificial Intelligence and Statistics*, pp. 805–813. PMLR, 2012.
- Ronald Ortner, Pratik Gajane, and Peter Auer. Variational regret bounds for reinforcement learning. In *Uncertainty in Artificial Intelligence*, pp. 81–90. PMLR, 2020.
- Ian Osband and Benjamin Van Roy. On lower bounds for regret in reinforcement learning. *arXiv preprint arXiv:1608.02732*, 2016.
- Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommannan, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks. *arXiv preprint arXiv:1712.03632*, 2017.
- Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial reinforcement learning. In *International Conference on Machine Learning*, pp. 2817–2826. PMLR, 2017.
- Aviv Rosenberg and Yishay Mansour. Online convex optimization in adversarial markov decision processes. In *International Conference on Machine Learning*, pp. 5478–5486. PMLR, 2019.
- Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Yoan Russac, Claire Vernade, and Olivier Cappé. Weighted linear bandits for non-stationary environments. In *Advances in Neural Information Processing Systems*, pp. 12017–12026, 2019.
- Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Lior Shani, Yonathan Efroni, and Shie Mannor. Adaptive trust region policy optimization: Global convergence and faster rates for regularized mdps. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 5668–5675, 2020.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Ahmed Touati and Pascal Vincent. Efficient learning in non-stationary linear markov decision processes. *arXiv preprint arXiv:2010.12870*, 2020.
- Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019a.
- Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Provably efficient reinforcement learning with general value function approximation. *arXiv preprint arXiv:2005.10804*, 2020.
- Yining Wang, Ruosong Wang, Simon S Du, and Akshay Krishnamurthy. Optimism in reinforcement learning with generalized linear function approximation. *arXiv preprint arXiv:1912.04136*, 2019b.
- Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019a.
- Lin F Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. *arXiv preprint arXiv:1902.04779*, 2019b.
- Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. On the global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. *arXiv preprint arXiv:1907.06246*, 2019.
- Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity. In *18th Annual Symposium on Foundations of Computer Science (sfcs 1977)*, pp. 222–227. IEEE, 1977.
- Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pp. 10978–10989. PMLR, 2020.
- Peng Zhao and Lijun Zhang. Non-stationary linear bandits revisited. *arXiv preprint arXiv:2103.05324*, 2021.
- Peng Zhao, Lijun Zhang, Yuan Jiang, and Zhi-Hua Zhou. A simple approach for non-stationary linear bandits. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics, AISTATS*, volume 2020, 2020.
- Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted MDPs with feature mapping. *arXiv preprint arXiv:2006.13165*, 2020a.
- Huozhi Zhou, Jinglin Chen, Lav R Varshney, and Ashish Jagmohan. Nonstationary reinforcement learning with linear function approximation. *arXiv preprint arXiv:2010.04244*, 2020b.
- Alexander Zimin and Gergely Neu. Online learning in episodic markovian decision processes by relative entropy policy search. In *Neural Information Processing Systems 26*, 2013.