
Online Sparse Reinforcement Learning

Botao Hao
Deepmind

Tor Lattimore
Deepmind

Csaba Szepesvári
University of Alberta/Deepmind

Mengdi Wang
Princeton University

Abstract

We investigate the hardness of online reinforcement learning in fixed horizon, sparse linear Markov decision process (MDP), with a special focus on the high-dimensional regime where the ambient dimension is larger than the number of episodes. Our contribution is two-fold. First, we provide a lower bound showing that linear regret is generally unavoidable in this case, even if there exists a policy that collects well-conditioned data. The lower bound construction uses an MDP with a fixed number of states while the number of actions scales with the ambient dimension. Note that when the horizon is fixed to one, the case of linear stochastic bandits, the linear regret can be avoided. Second, we show that if the learner has oracle access to a policy that collects well-conditioned data then a variant of Lasso fitted Q-iteration enjoys a nearly dimension free regret of $\tilde{O}(s^{2/3}N^{2/3})$ where N is the number of episodes and s is the sparsity level. This shows that in the large-action setting, the difficulty of learning can be attributed to the difficulty of finding a good exploratory policy.

1 INTRODUCTION

Sparse models in classical statistics often yield the best of both worlds: high representation power is achieved by including many features while sparsity leads to efficient estimation. There is a growing interest in applying the tools developed by statisticians to sequential settings such as contextual bandits and reinforcement learning (RL). As we now explore, in online RL this leads to a number of deli-

cate trade-offs between assumptions and sample complexity. The use of sparsity in reinforcement learning (RL) has been explored before in the context of policy evaluation or policy optimization in the batch setting [Kolter and Ng, 2009, Geist and Scherrer, 2011, Hoffman et al., 2011, Painter-Wakefield and Parr, 2012, Ghavamzadeh et al., 2011, Hao et al., 2020a]. As far as we know, there has been very little work on the role of sparsity in online RL. In batch RL, the dataset is given a priori and the focus is typically on evaluating a given target policy or learning a near-optimal policy. By contrast, the central question in online RL is how to sequentially interact with the environment to balance the trade-off between exploration and exploitation, measured here by the cumulative regret. We ask the following question:

Under what circumstances does sparsity help when minimizing regret in online RL?

In sparse linear regression, the optimal estimation error rate generally scales polynomially with the sparsity s and only logarithmically in the ambient dimension d [Wainwright, 2019]. This is guaranteed under the sufficient and almost necessary condition that the data covariance matrix is well-conditioned, usually referred to restricted eigenvalue condition [Bickel et al., 2009] or compatibility condition [Van De Geer et al., 2009].

The ‘almost necessary’ nature of the conditions for efficient estimation with sparsity leads to an unpleasant situation when minimizing regret. Even in sparse linear bandits, the worst-case regret is known to depend polynomially on the ambient dimension [Lattimore and Szepesvári, 2020, §24.3]. The reason is simple. By definition, a learner with small regret must play mostly the optimal action, which automatically leads to poorly conditioned data. Hence, making the right assumptions is essential in the high-dimensional regime where d is large relative to the time horizon. A number of authors have considered the contextual setting, where the regret can be made dimension free by making judicious assumptions on the context distribution [Bastani and Bayati, 2020, Wang et al., 2018, Kim and Paik, 2019,

Ren and Zhou, 2020, Wang et al., 2020].

When lifting assumptions from the bandit literature to RL it is essential to ensure that (a) the assumptions still help and (b) the assumptions remain reasonable. In some sense, our lower bound shows that a typical assumption that helps in linear bandits is by itself insufficient in RL. Specifically, in linear bandits the existence of a policy that collects well-conditioned data is sufficient for dimension free regret. In RL this is no longer true because finding this policy may not be possible without first learning the transition structure, which cannot be done efficiently without well-conditioned data, which yields an irresolvable chicken-and-egg problem.

Contribution We study online RL in episodic linear MDPs with ambient dimension d , sparsity s , episode length H and number of episodes N . Our contribution is two-fold:

- Our first result is a lower bound showing that $\Omega(Hd)$ regret is unavoidable in the worse-case when the dimension is large, even if the MDP transition kernel can be exactly represented by a sparse linear model and there exists an exploratory policy that collects well-conditioned data. The technical contribution is to craft a new class of hard-to-learn episodic MDPs. To overcome the difficulties caused by deterministic transitions from the constructed MDPs, we develop a novel stopping-time argument when calculating the KL-divergence.
- Our second result shows that if the learner has oracle access to an exploratory policy that collects well-conditioned data, then online Lasso fitted-Q-iteration in combination with the explore-then-commit template achieves a regret upper bound of $\tilde{O}(H^{4/3}s^{2/3}N^{2/3})$. The proof requires a non-trivial extension of high-dimensional statistics to Markov dependent data. As far as we know, this is the first regret bound that has no polynomial dependency on the feature dimension d in online RL.

1.1 Related work

Regret guarantees for online RL have received considerable attention in recent years. In episodic tabular MDPs with a homogeneous transition kernel, Azar et al. [2017] proved a minimax optimal regret of $O(\sqrt{H^2|S||A|N})$ achieved by a model-based algorithm. Jin et al. [2018] showed an $O(\sqrt{H^4|S||A|N})$ regret bound for Q-learning with inhomogeneous transition kernel. Under a linear MDP assumption, Jin et al. [2019] showed an $O(\sqrt{d^3H^4N})$ regret bound for an optimistic version of least-squares value iteration.

Under a linear kernel MDP assumption [Zhou et al., 2020], Yang and Wang [2020] obtained an $O(dH^{5/2}\sqrt{N})$ regret bound by a model-based algorithm while Cai et al. [2019] obtained an $O(dH^2\sqrt{N})$ regret bound using an optimistic version of least-squares policy iteration. Zanette et al. [2020] derived an $O(d^2H^{5/2}\sqrt{N})$ regret bound for randomized least-squares value iteration. None of these works considered sparsity, and consequentially the aforementioned regret bounds all have polynomial dependency on d .

Jiang et al. [2017] and Sun et al. [2019] design algorithms for learning in RL problems with low Bellman/Witness rank, which includes sparse linear RL as a special case and obtain $O(\text{poly}(s, A, H, \log(d)))$ sample complexity where A is the number of actions. More recently, FLAMBE [Agarwal et al., 2020] achieves $O(\text{poly}(s, A, H, \log(d)))$ sample complexity in a low-rank MDP setting. It is worth mentioning that although the above results have no polynomial dependency on d , the sample complexity unavoidably involves polynomial dependency on the number of actions.

There are several previous works focusing on sparse linear/contextual bandits that can be viewed as a simplified online RL problem. Abbasi-Yadkori et al. [2012] proposed an online-to-confidence-set conversion approach that achieves an $O(\sqrt{sdN})$ regret upper bound, where s is a known upper bound on the sparsity. The algorithm is not computationally-efficient, a deficit that is widely believed to be unavoidable. A matching lower bound is also known, which means polynomial dependence on d is generally unavoidable without additional assumptions [Lattimore and Szepesvári, 2020, §24.3]. More recently, under the condition that the feature vectors admit a well-conditioned exploration distribution, Hao et al. [2020b] proves a dimension-free $\Omega(s^{1/3}N^{2/3})$ regret lower bound in the high-dimensional regime that can be matched by an explore-then-commit algorithm. In the contextual setting, where the action set changes from round to round, several works imposed various of careful assumptions on the context distribution such that polynomial dependency on d can be removed [Bastani and Bayati, 2020, Wang et al., 2018, Kim and Paik, 2019, Ren and Zhou, 2020, Wang et al., 2020]. As far as we can tell, however, these assumptions are not easily extended to the MDP setting, where the contextual information available to the learner is not independent and identically distributed.

The use of feature selection in offline RL has also been investigated in a number of prior works. Kolter and Ng [2009], Geist and Scherrer [2011], Hoffman et al. [2011], Painter-Wakefield and Parr [2012], Liu et al. [2012] studied on-policy/off-policy evaluation with ℓ_1 -regularization for temporal-difference (TD) learning. Ghavamzadeh et al.

[2011] and Geist et al. [2012] proposed Lasso-TD to estimate the value function in Markov reward processes and derived finite-sample statistical analysis. However, the aforementioned results can not be extended to online setting directly. Hao et al. [2020a] provides nearly optimal statistical analysis for sparse off-policy evaluation/optimization. One exception by Ibrahimi et al. [2012], who derived an $O(p\sqrt{N})$ regret bound in high-dimensional sparse linear quadratic systems where p is the dimension of the state space.

2 PRELIMINARY

Notation. Denote by $\sigma_{\min}(X)$ and $\sigma_{\max}(X)$ the smallest and largest eigenvalues of a symmetric matrix X . Let $[n] = \{1, 2, \dots, n\}$. The relations \lesssim and \gtrsim stand for “approximately less/greater than” and are used to omit constant and poly-logarithmic factors. We use $\tilde{O}(\cdot)$ to omit polylog factors. For a finite set \mathcal{S} , let $\Delta_{\mathcal{S}}$ be the set of probability distributions over \mathcal{S} .

2.1 Problem definition

Episodic MDP. A finite episodic Markov decision process (MDP) is a tuple $(\mathcal{X}, \mathcal{A}, H, P, r)$ with \mathcal{X} the state-space, \mathcal{A} the action space, H the episode length, $P : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_{\mathcal{X}}$ the transition kernel and $r : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ the reward function. As is standard, we assume that \mathcal{X} and \mathcal{A} are finite and that the reward function is known. We write $P(x'|x, a)$ for the probability of transitioning to state x' when taking action a in state x . A learner interacts with an episodic MDP as follows. In each episode, an initial state x_1 is sampled from an initial distribution $\xi_0 \in \Delta_{\mathcal{X}}$. Then, in each step $h \in [H]$, the learner observes a state $x_h \in \mathcal{X}$, takes an action $a_h \in \mathcal{A}$, and receives a deterministic reward $r(x_h, a_h)$. Then, the system evolves to a random next state x_{h+1} according to distribution $P(\cdot|x_h, a_h)$. The episode terminates when x_{H+1} is reached.

We define a (stationary) policy as a function $\pi : \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$, that maps states to distributions over actions. A nonstationary policy is a sequence of maps from histories to probability distributions over actions. For each $h \in [H]$ and policy π , the value function $V_h^\pi : \mathcal{X} \rightarrow \mathbb{R}$ is defined as the expected value of cumulative rewards received under policy π when starting from an arbitrary state at h th step; that is,

$$V_h^\pi(x) := \mathbb{E}^\pi \left[\sum_{h'=h}^H r(x_{h'}, a_{h'}) \middle| x_h = x \right],$$

where $a_{h'} \sim \pi(\cdot|x_{h'})$, $x_{h'+1} \sim P(\cdot|x_{h'}, a_{h'})$, and \mathbb{E}^π denotes the expectation over the sample path generated under

policy π . Accordingly, we also define the action-value function $Q_h^\pi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ which gives the expected cumulative reward when the learner starts from an arbitrary state-action pair at the h th step and follows policy π afterwards:

$$Q_h^\pi(x, a) := r(x, a) + \mathbb{E}^\pi \left[\sum_{h'=h+1}^H r(x_{h'}, a_{h'}) \middle| x_h = x, a_h = a \right].$$

Note, the conditioning in the above definitions is not quite innocent. In this form the value function is not well defined for states x that are not reachable by a given policy. This is easily rectified by defining the value function in terms of the Bellman equation or by being more rigorous about the probability space. The above definitions are standard in the literature and are left as is for reader’s convenience.

Bellman equation. Since the action space and episode length are both finite, there always exists an optimal policy π^* which gives the optimal value $V_h^*(x) = \sup_{\pi} V_h^\pi(x)$ for all $x \in \mathcal{X}$ and $h \in [H]$ [Puterman, 2014, Szepesvári, 2010]. We denote the Bellman operator as

$$[\mathcal{T}V](x, a) := r(x, a) + \mathbb{E}_{x' \sim P(\cdot|x, a)}[V(x')],$$

and the Bellman equation for policy π becomes

$$\begin{aligned} Q_h^\pi(x, a) &= [\mathcal{T}V_{h+1}^\pi](x, a), \\ V_h^\pi(x) &= \mathbb{E}_{a \sim \pi(\cdot|x)}[Q_h^\pi(x, a)], V_{H+1}^\pi(x) = 0, \end{aligned} \quad (2.1)$$

which holds for all $(x, a) \in \mathcal{X} \times \mathcal{A}$. Similarly, the Bellman optimality equation is

$$\begin{aligned} Q_h^*(x, a) &= [\mathcal{T}V_{h+1}^*](x, a), \\ V_h^*(x) &= \max_{a \in \mathcal{A}} Q_h^*(x, a), V_{H+1}^*(x) = 0. \end{aligned} \quad (2.2)$$

Cumulative regret. In the online setting, the learner aims to minimize the cumulative regret by interacting with the environment over a number of episodes. At the beginning of the n th episode, an initial state x_1^n is sampled from ξ_0 and the agent executes policy π_n . We measure the performance of the learner over N episodes by the cumulative regret:

$$R_N = \sum_{n=1}^N (V_1^*(x_1^n) - V_1^{\pi_n}(x_1^n)). \quad (2.3)$$

The cumulative regret measures the expected loss of following the policy produced by the learner instead of the optimal policy. Therefore, the learner aims to follow a sequence of policies π_1, \dots, π_N such that the cumulative regret is minimized.

2.2 Sparse linear MDPs

Before we introduce sparse linear MDPs, we need to settle on a definition of a linear MDP. Let $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a feature map which assigns to each state-action pair a d -dimensional feature vector. A feature map combined with a parameter vector $w \in \mathbb{R}^d$ gives rise to the linear function $g_w : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ defined by $g_w(x, a) = \phi(x, a)^\top w$ and the subspace $\mathcal{G}_\phi = \{g_w : w \in \mathbb{R}^d\} \subseteq \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$. Given a policy π and function $f : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$, let $\tilde{\mathcal{T}}_\pi f : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ be the function defined by

$$[\tilde{\mathcal{T}}_\pi f](x, a) = r(x, a) + \mathbb{E}_{x' \sim P(\cdot|x, a), a \sim \pi(a|x')} [f(x', a)].$$

We call an MDP *linear* if \mathcal{G}_ϕ is closed under $\tilde{\mathcal{T}}_\pi$ for all policies π .¹ Yang and Wang [2019] and Jin et al. [2019] have shown that this is equivalent to assuming $P(x'|x, a) = \sum_{k \in [d]} \phi_k(x, a) \psi_k(x')$, for some functions $\psi_1, \dots, \psi_d : \mathcal{X} \rightarrow \mathbb{R}$ and all pairs of (x, a) . Note, the feature map ϕ is always assumed to be known to the learner. As far as we know, this notion of linearity was introduced by Bellman et al. [1963], Schweitzer and Seidmann [1985], who were motivated by the problem of efficiently computing the optimal policy for a known MDP with a large state-space.

When little priori information is available on how to choose the features, agnostic choices often lead to dimensions which can be as large as the number of episodes N . Without further assumptions, no procedure can achieve nontrivial performance guarantees, even when just considering simple prediction problems (e.g., predicting immediate rewards). However, effective learning with many more features than the sample-size is possible when only $s \ll d$ features are relevant. This motivates our definition of a sparse linear MDP.

Definition 2.1 (Sparse linear MDP). Fix a feature map $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and assume the episodic MDP \mathcal{M} is linear in ϕ . We say \mathcal{M} is (s, ϕ) -sparse if there exists an active set $\mathcal{K} \subseteq [d]$ with $|\mathcal{K}| \leq s$ and some functions $\psi(\cdot) = (\psi_k(\cdot))_{k \in \mathcal{K}}$ such that for all pairs of $(x, a) : P(x'|x, a) = \sum_{k \in \mathcal{K}} \phi_k(x, a) \psi_k(x')$.

3 HARDNESS OF ONLINE SPARSE RL

In this section we illustrate the fundamental hardness of online sparse RL in the high-dimensional regime by proving a minimax regret lower bound. The high-dimensional

¹A different definition is called linear kernel MDP that the MDP transition kernel can be parameterized by a small number of parameters [Yang and Wang, 2020, Cai et al., 2019, Zanette et al., 2020, Zhou et al., 2020]

regime is referred to $N \leq d$. We first introduce a notion of an exploratory policy.

Definition 3.1 (Exploratory policy). Let Σ^π be the expected uncentered covariance matrix induced by policy π and feature map ϕ , given by

$$\Sigma^\pi := \mathbb{E}^\pi \left[\frac{1}{H} \sum_{h=1}^H \phi(x_h, a_h) \phi(x_h, a_h)^\top \right], \quad (3.1)$$

where $x_1 \sim \xi_0, a_h \sim \pi(\cdot|x_h), x_{h+1} \sim P(\cdot|x_h, a_h)$ and \mathbb{E}^π denotes expectation over the sample path generated under policy π . We call a policy π *exploratory* if $\sigma_{\min}(\Sigma^\pi) > 0$.

Remark 3.2. In the tabular case, we can choose $\phi(x, a)$ as a basis vector in $\mathbb{R}^{|\mathcal{X}| \times |\mathcal{A}|}$. Let $\mu^\pi(x, a)$ be the frequency of visitation for state-action pair (x, a) under policy π and initial distribution ξ_0 :

$$\mu^\pi(x, a) = \frac{1}{H} \sum_{h=1}^H \mathbb{E}^\pi [\mathbb{I}((x_h, a_h) = (x, a))].$$

Then $\sigma_{\min}(\Sigma^\pi) > 0$ implies $\min_{x,a} \mu^\pi(x, a) > 0$. This means an exploratory policy in the tabular case will have positive visitation probability of each state-action pair.

The next theorem is a kind of minimax lower bound for online sparse RL. The key steps of the proof follow, with details and technical lemmas deferred to the appendix.

Theorem 3.3 (Minimax lower bound in high-dimensional regime). For any algorithm π , there exists a sparse linear MDP \mathcal{M} and associated exploratory policy π_e for which $\sigma_{\min}(\Sigma^{\pi_e})$ is a strictly positive universal constant independent of N and d , such that for any $N \leq d$,

$$R_N \geq \frac{1}{128} H d.$$

This theorem states that even if the MDP transition kernel can be exactly represented by a sparse linear model and there exists an exploratory policy, the learner could still suffer linear regret in the high-dimensional regime. This is in stark contrast to linear bandits, where the existence of an exploratory policy is sufficient for dimension-free regret. The problem in RL is that *finding* the exploratory policy can be very hard.

Remark 3.4. Our lower bound construction may suggest some general principles where we can exploit the sparsity easily: what seems to be essential for exploiting sparsity is that the learner can collect well conditioned data without needing to explore the whole MDP, i.e., some prior knowledge is needed, or the class of MDPs must be restricted to those where many actions provide well conditioned data.

Proof of Theorem 3.3. The proof uses the standard information theoretic machinery, but with a novel hard-to-learn MDP construction and KL divergence calculation based on a stopping time argument. The intuition is to construct an informative state with only one of a large set of actions leading to the informative state deterministically. And the exploratory policy has to visit that informative state to produce well-conditioned data. In order to find this informative state, the learner should take a large number of trials that will suffer high regret.

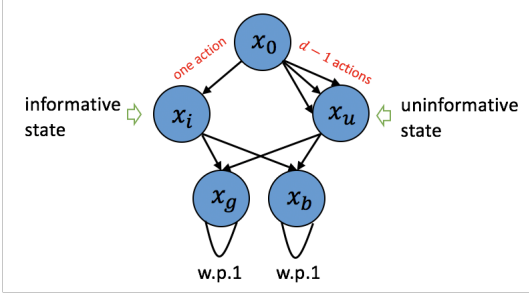


Figure 1: A hard-to-learn MDP instance that includes an informative state and an uninformative state.

Step 1: Construct a set of hard MDP instances. Let the state space \mathcal{X} consists of $\{x_0, x_i, x_u, x_g, x_b\}$. Here, x_0 is the initial state, x_i and x_u refer to informative and uninformative states, x_g and x_b refer to high-reward and low-reward states. Construct d different hard MDP instances: $\{\mathcal{M}_1, \dots, \mathcal{M}_d\}$ and they only differ at which action brings the learner to x_i . For each MDP $\mathcal{M}_k, k \in [d]$

$$\theta = \left(\underbrace{\varepsilon, \dots, \varepsilon}_{s-1}, 0, \dots, 0, \frac{1}{2} \right)^\top \in \mathbb{R}^d, \quad (3.2)$$

where $\varepsilon > 0$ is a small constant to be tuned later, and $\bar{\theta}^{(k)} \in \mathbb{R}^{2d+2}$ as

$$\bar{\theta}^{(k)} = \left(\theta^\top, 1, 1, \underbrace{0, \dots, 0}_{k-1}, 1, \underbrace{0, \dots, 0}_{d-k} \right)^\top. \quad (3.3)$$

We specify the transition probability of \mathcal{M}_k in the following steps:

1. Let the initial state x_0 associated with d actions as $\mathcal{A}_1 = \{a_1^0, \dots, a_d^0\}$. The transitions from x_0 to either x_i or x_u are *deterministic*. In MDP \mathcal{M}_k , only taking action a_k^0 brings the learner to x_i , and taking any other action *except* a_k^0 brings the learner to x_u . This information is parameterized into the last d coordinates of $\bar{\theta}^{(k)}$. We defer the detailed construction of features to Appendix B.1.

2. We construct a feature set \mathcal{S} associated to x_u and a feature set \mathcal{H} associated to x_i :

$$\mathcal{S} = \left\{ z \in \mathbb{R}^d \mid z_d = 0, z_j \in \{-1, 0, 1\} \right. \\ \left. \text{for } j \in [d-1], \|z\|_1 = s-1 \right\}, \\ \mathcal{H} = \left\{ z \in \mathbb{R}^d \mid z_j \in \{-1, 1\} \text{ for } j \in [d-1], z_d = 1 \right\}.$$

Let $\mathcal{A}_2 = \{a_1^u, \dots, a_{|\mathcal{S}|}^u\}$ be the action set associated with x_u and $\mathcal{A}_3 = \{a_1^i, \dots, a_{|\mathcal{H}|}^i\}$ be the action set associated with x_i . We write $\varphi(x_u, a_j^u)$ as the j th element in \mathcal{S} and $\varphi(x_i, a_j^i)$ as the j th element in \mathcal{H} . Denote

$$\phi(x_i, a_j^i) = \left(\varphi(x_i, a_j^i)^\top, \underbrace{0, \dots, 0}_{d+2}, 1 \right)^\top \in \mathbb{R}^{2d+3}$$

and $\psi(x_b) = (\bar{\theta}^{(k)\top}, 0) \in \mathbb{R}^{2d+3}$, $\psi(x_g) = (-\bar{\theta}^{(k)\top}, 1) \in \mathbb{R}^{2d+3}$. At informative state x_i , the learner can take action $a_j^i \in \mathcal{A}_3$ and transits to either x_g or x_b according to

$$P(x_g | x_i, a_j^i) = \varphi(x_i, a_j^i)^\top \theta = \phi(x_i, a_j^i)^\top \psi(x_g), \\ P(x_b | x_i, a_j^i) = 1 - \varphi(x_i, a_j^i)^\top \theta = \phi(x_i, a_j^i)^\top \psi(x_b),$$

that satisfy the sparse linear MDP assumption. We specify ϕ, ψ similarly when the learner at x_u .

3. At x_g or x_b , the learner will stay the current state for the rest of current episode no matter what actions to take.

One can verify that the above construction so far satisfies the sparse linear MDP assumption in Definition 2.1. In the end, the reward function is set to be $r(x, a) = 1$ if $x = x_g$ and $r(x, a) = 0$ otherwise. We now finish the construction of all the essential ingredients of $\{\mathcal{M}_1, \dots, \mathcal{M}_d\}$.

Remark 3.5. For \mathcal{M}_k , the overall action set will be $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \mathcal{A}_3$. Now we specify the transitions that we have not mentioned so far. At x_0 , all the actions from \mathcal{A}_2 and \mathcal{A}_3 bring the learner to x_u . At x_i , all the actions from \mathcal{A}_1 and \mathcal{A}_2 bring the learner to either x_g or x_b with the same probability as a_1^i . At x_u , actions from \mathcal{A}_1 and \mathcal{A}_3 bring the learner to either x_g or x_b with the same probability as a_1^u .

Step 2: Construct an alternative set of MDPs. For each $k \in [d]$, the second step is to construct an alternative MDP $\tilde{\mathcal{M}}_k$ that is hard to distinguish from \mathcal{M}_k and for which the optimal policy for \mathcal{M}_k is suboptimal for $\tilde{\mathcal{M}}_k$ and vice versa. Fix a sequence of policies $\{\pi_1, \dots, \pi_N\}$. Let $\mathcal{D}_n = (S_1^n, A_1^n, \dots, S_H^n, A_H^n)$ be the sequence of state-action pairs in n th episode produced by π_n . Define $\mathcal{F}_h^n = \sigma(\mathcal{D}_1, \dots, \mathcal{D}_{h-1}, S_1^n, A_1^n, \dots, S_{h-1}^n, A_{h-1}^n, S_h^n)$. Let $\mathbb{F} =$

$(\mathcal{F}_h^n)_{h \in [H], n \in [N]}$ be a filtration. Define the stopping time with respect to \mathbb{F} : $\tau_k = N \wedge \min\{n : A_1^n = a_k^0\}$ that is the first episode the learner reaches the informative state. In other words, for $n \leq \tau_k - 1$, the learner always transits to x_u from x_0 . At x_u , the learner acts similarly as facing linear bandits where the number of arms is $|\mathcal{S}|$.

For $k \in [d]$, let $\mathbb{P}_k, \widetilde{\mathbb{P}}_k$ be the laws of $\mathcal{D}_1, \dots, \mathcal{D}_{\tau_k-1}$ induced by the interaction of $\{\pi_1, \dots, \pi_N\}$ and $\mathcal{M}_k, \widetilde{\mathcal{M}}_k$ accordingly. Let $\mathbb{E}_k, \widetilde{\mathbb{E}}_k$ be the corresponding expectation operators. In addition, denote a set \mathcal{S}' as

$$\begin{aligned} \mathcal{S}' = & \left\{ z \in \mathbb{R}^d \mid \|z\|_1 = s - 1, \right. \\ & z_j \in \{-1, 0, 1\} \text{ for } j \in \{s, s+1, \dots, d-1\}, \\ & \left. z_j = 0 \text{ for } j \in \{1, \dots, s-1, d\} \right\}. \end{aligned} \quad (3.4)$$

Then we let

$$\widetilde{z}^{(k)} = \operatorname{argmin}_{z \in \mathcal{S}'} \mathbb{E}_k \left[\sum_{n=1}^{\tau_k} \langle \varphi(S_2^n, A_2^n), z \rangle^2 \right], \quad (3.5)$$

and construct the alternative $\widetilde{\theta}^{(k)} = \theta + 2\varepsilon \widetilde{z}^{(k)}$ where ε appears in Eq. (3.2). This is in contrast with $\bar{\theta}^{(k)}$ in Eq. (3.3) that specifies the original MDP \mathcal{M}_k . All the other ingredients of $\widetilde{\mathcal{M}}_k$ are the same as \mathcal{M}_k . Thus, we have constructed an alternative set of MDPs.

Step 3: Regret decomposition. Let $R_N(\mathcal{M}_k)$ be the cumulative regret of a sequence of policies $\{\pi_1, \dots, \pi_N\}$ interacting with MDP \mathcal{M}_k for N episodes. Recall that from the definition in Eq. (2.3), we have

$$R_N(\mathcal{M}_k) = \sum_{n=1}^N \left(V_1^{\pi^*}(x_1^n) - V_1^{\pi_n}(x_1^n) \right).$$

Denote $a^* = \operatorname{argmax}_{a_j^u \in \mathcal{A}_2} \varphi(x_u, a_j^u)^\top \theta$ be the optimal action when the learner is at x_u . The optimal policy π^* of MDP \mathcal{M}_k behaves in the following way for each episode:

- At state x_0 , the optimal policy takes an arbitrary action except a_k^0 to state x_u . There is no reward collected so far.
- At state x_u , the optimal policy takes action a^* and transits to good state x_g with probability $\varphi(x_u, a^*)^\top \theta$ or bad state x_b with probability $1 - \varphi(x_u, a^*)^\top \theta$.
- The learner stays at the current state for the rest of current episode.

Then the value function of π^* at n th episode is

$$\begin{aligned} V_1^{\pi^*}(x_1^n) &= (H-1)\mathbb{P}(A_2^n = a^*) \\ &= (H-1)\varphi(x_u, a^*)^\top \theta = (H-1)(s-1)\varepsilon. \end{aligned}$$

We decompose $R_N(\mathcal{M}_k)$ according to the stopping time τ_k :

$$\begin{aligned} R_N(\mathcal{M}_k) &\geq \sum_{n=1}^{\tau_k-1} \left(V_1^{\pi^*}(x_1^n) - V_1^{\pi_n}(x_1^n) \right) \\ &\geq \frac{H}{8} \mathbb{E}_k \left[\tau_k s \varepsilon - \sum_{n=1}^{\tau_k-1} \langle \varphi(S_2^n, A_2^n), \theta \rangle \right] \\ &= \frac{H}{8} \mathbb{E}_k \left[\tau_k s \varepsilon - \sum_{n=1}^{\tau_k-1} \sum_{j=1}^{s-1} \varphi_j(x_u, A_2^n) \varepsilon \right], \end{aligned} \quad (3.6)$$

where the last equation is due to S_2^n is always x_u until the stopping time τ_k . Define an event

$$\mathcal{D}_k = \left\{ \sum_{n=1}^{\tau_k-1} \sum_{j=1}^{s-1} \varphi_j(x_u, A_2^n) \leq \frac{\tau_k s}{2} \right\}.$$

The next claim shows that when \mathcal{D}_k occurs, the regret is large in MDP \mathcal{M}_k , while if it does not occur, then the regret is large in MDP $\widetilde{\mathcal{M}}_k$. The detailed proof is deferred to Appendix B.2.

Claim 3.6. Regret lower bounds with respect to event \mathcal{D}_k :

$$\begin{aligned} R_N(\mathcal{M}_k) + R_N(\widetilde{\mathcal{M}}_k) &\geq \frac{Hs\varepsilon}{8} \left(\mathbb{E}_k[\tau_k] \right. \\ &\quad \left. + \widetilde{\mathbb{E}}_k[\tau_k \mathbb{I}(\mathcal{D}_k^c)] - \mathbb{E}_k[\tau_k \mathbb{I}(\mathcal{D}_k^c)] \right). \end{aligned}$$

We construct an additional MDP \mathcal{M}_0 such that when the learner is at x_0 , no matter what actions to take, the learner will always transit to the uninformative state x_u . All the other structures remain the same with $\{\mathcal{M}_1, \dots, \mathcal{M}_d\}$. Let \mathbb{P}_0 be the laws of $\mathcal{D}_1, \dots, \mathcal{D}_{\tau_k-1}$ induced by the interaction of π and \mathcal{M}_0 and let \mathbb{E}_0 be the corresponding expectation operators. Then from Pinsker's inequality (Lemma C.5 in the Appendix), for any $k \in [d]$,

$$\begin{aligned} \left| \mathbb{E}_0[\tau_k] - \mathbb{E}_k[\tau_k] \right| &\leq N \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}_0 \| \mathbb{P}_k)}, \\ \left| \widetilde{\mathbb{E}}_k[\tau_k \mathbb{I}(\mathcal{D}_k^c)] - \mathbb{E}_k[\tau_k \mathbb{I}(\mathcal{D}_k^c)] \right| &\leq N \sqrt{\frac{1}{2} \text{KL}(\widetilde{\mathbb{P}}_k \| \mathbb{P}_k)}, \end{aligned}$$

where $\text{KL}(\mathbb{P}, \mathbb{P}')$ is the KL divergence between probability measures \mathbb{P} and \mathbb{P}' . Combining with Claim 3.6, we have

$$\begin{aligned} R_N(\mathcal{M}_k) + R_N(\widetilde{\mathcal{M}}_k) &\geq \frac{Hs\varepsilon}{8} \left(\mathbb{E}_0[\tau_k] \right. \\ &\quad \left. - d \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}_0 \| \mathbb{P}_k)} - d \sqrt{\frac{1}{2} \text{KL}(\widetilde{\mathbb{P}}_k \| \mathbb{P}_k)} \right), \end{aligned} \quad (3.7)$$

where we consider the high-dimensional regime such that $N \leq d$.

Step 4: Calculating the KL divergence. We make use of the following bound on the KL divergence between $\tilde{\mathbb{P}}_k$ and $\mathbb{P}_k, \mathbb{P}_0$ and \mathbb{P}_k , which formalises the intuitive notion of information. When the KL divergence is small, the algorithm is unable to distinguish the two environments. The detailed proof is deferred to Appendix B.3.

Claim 3.7. The KL divergences between $\tilde{\mathbb{P}}_k$ and $\mathbb{P}_k, \mathbb{P}_0$ and \mathbb{P}_k are upper bounded by the following when $N \leq d$:

$$\text{KL}(\tilde{\mathbb{P}}_k \| \mathbb{P}_k) \leq 8\varepsilon^2(s-1)^2, \text{KL}(\mathbb{P}_0 \| \mathbb{P}_k) = 0. \quad (3.8)$$

Combining with Eq. (B.11) and summing over the set of MDPs $\{\mathcal{M}_1, \dots, \mathcal{M}_d\}$,

$$\sum_{k=1}^d \left(R_N(\mathcal{M}_k) + R_N(\tilde{\mathcal{M}}_k) \right) \geq \frac{Hs\varepsilon}{8} \left(\sum_{k=1}^d \mathbb{E}_0[\tau_k] - d^2 \sqrt{8\varepsilon^2 s^2} \right).$$

Picking $\varepsilon = 1/(8s)$, we have

$$\sum_{k=1}^d \left(R_N(\mathcal{M}_k) + R_N(\tilde{\mathcal{M}}_k) \right) \geq \frac{H}{32} \left(\sum_{k=1}^d \mathbb{E}_0[\tau_k] - \frac{d^2}{4} \right).$$

Step 5: Summary. From the definition of the stopping time, one can see $\sum_{k=1}^d \mathbb{E}_0(\tau_k) \geq \sum_{k=1}^d k \geq d^2/2$. Therefore,

$$\sum_{k=1}^d \left(R_N(\mathcal{M}_k) + R_N(\tilde{\mathcal{M}}_k) \right) \geq \frac{1}{128} H d^2.$$

Among two sets of MDPs $\{\mathcal{M}_k\}_{k=1}^d$ and $\{\tilde{\mathcal{M}}_k\}_{k=1}^d$, for any sequence of policies $\{\pi_1, \dots, \pi_N\}$, there must exist a MDP \mathcal{M}_k such that

$$R_N(\mathcal{M}_k) \geq \frac{1}{128} H d.$$

This finishes the proof. \square

4 ONLINE LASSO FITTED-Q-ITERATION

In this section we prove that if the learner has oracle access to an exploratory policy, the online Lasso fitted-Q-iteration (Lasso-FQI) algorithm can have a dimension-free $\tilde{O}(N^{2/3})$ regret upper bound.

We first introduce the online Lasso-FQI. Suppose the learner has the oracle access to an exploratory policy π_e (defined in Definition 3.1). The algorithm uses the explore-then-commit template and includes the following three phases:

- **Exploration phase.** The exploration phase includes N_1 episodes where N_1 will be chosen later based on

regret bound and can be factorized as $N_1 = RH$, where $R > 1$ is an integer. At the beginning of each episode, the agent receives an initial state drawn from ξ_0 and executes the rest steps following the exploratory policy π_e . Let the dataset collected in the exploration stage as \mathcal{D} .

- **Learning phase.** Split \mathcal{D} into H folds: $\{\mathcal{D}_1, \dots, \mathcal{D}_H\}$ and each fold consists of R episodes. Based on the exploratory dataset \mathcal{D} , the agent executes an extension of fitted-Q-iteration [Ernst et al., 2005, Antos et al., 2008] combining with Lasso [Tibshirani, 1996] for feature selection. To define the algorithm, it is useful to introduce $Q_w(x, a) = \phi(x, a)^\top w$. For $a < b$, we also define the operator $\Pi_{[a,b]} : \mathbb{R} \rightarrow [a, b]$ that projects its input to $[a, b]$, i.e., $\Pi_{[a,b]}(x) = \max(\min(x, b), a)$. Initialize $\hat{w}_{H+1} = 0$. At each step $h \in [H]$, we fit \hat{w}_h through Lasso:

$$\hat{w}_h = \underset{w}{\operatorname{argmin}} \lambda_1 \|w\|_1 + \frac{1}{|\mathcal{D}_h|} \sum_{(x_i, a_i, x'_i) \in \mathcal{D}_h} (y_i - \phi(x_i, a_i)^\top w)^2, \quad (4.1)$$

where $y_i = \Pi_{[0,H]} \max_{a \in \mathcal{A}} Q_{\hat{w}_{h+1}}(x'_i, a)$ and λ_1 is a regularization parameter.

- **Exploitation phase.** For the rest $N - N_1$ episodes, the agent commits to the greedy policy with respect to the estimated Q-value $\{Q_{\hat{w}_h}\}_{h=1}^H$.

The full algorithm of online Lasso-FQI is summarized in Algorithm 1.

Remark 4.1. A key observation of Algorithm 1 is that the expected covariance matrix of data collected in the exploration phase could be well-conditioned due to the use of exploratory policy, e.g.,

$$\sigma_{\min} \left(\mathbb{E}^{\pi_e} \left[\frac{1}{N_1 H} \sum_{n=1}^{N_1} \sum_{h=1}^H \phi(x_h^n, a_h^n) \phi(x_h^n, a_h^n)^\top \right] \right) > 0.$$

This is the key condition to ensure the success of fast sparse feature selection in the learning and exploitation phases and eliminate the polynomial dependency of d in the cumulative regret.

Next we derive the regret guarantee for the online Lasso-FQI under the sparse linear MDP model. The definition of restricted eigenvalue $C_{\min}(\Sigma, s)$ and the proof is deferred to Appendix A.

Theorem 4.2 (Regret bound for online Lasso-FQI). Suppose the episodic MDP is (s, ϕ) -sparse as defined in

Algorithm 1 Online Lasso-FQI

-
- 1: **Input:** An episodic MDP $\mathcal{M} = (\mathcal{X}, \mathcal{A}, P, r, H)$, an exploratory policy π_e , exploration length N_1 , regularization parameter λ_1 ;
 - 2: **Initialize.** $\mathcal{D} = \emptyset$.
 - 3: **for** $n = 1, \dots, N_1$ **do**
 - 4: **for** $h = 1, \dots, H$ **do**
 - 5: Take the action $a_h^n = \pi_e(\cdot | x_h^n)$ and observe x_{h+1}^n .
 - 6: Let $\mathcal{D} = \mathcal{D} \cup \{x_h^n, a_h^n, x_{h+1}^n\}$.
 - 7: **end for**
 - 8: **end for**
 - 9: Partition the dataset \mathcal{D} into H folds such that each fold \mathcal{D}_h has R different episodes.
 - 10: Initialize $Q_{\hat{w}_{H+1}}(x, a) = 0$.
 - 11: **for** $h = H, \dots, 1$ **do**
 - 12: Calculate regression targets for each $(x_i, a_i, x'_i) \in \mathcal{D}_h$: $y_i = \Pi_{[0, H]} \max_{a \in \mathcal{A}} Q_{\hat{w}_{h+1}}(x'_i, a)$.
 - 13: Build training set $\{(x_i, a_i), y_i\}_{i \in \mathcal{D}_h}$ and fit \hat{w}_h through sparse linear regression in Eq. (4.1).
 - 14: **end for**
 - 15: **for** $n = N_1 + 1$ to N **do**
 - 16: **for** $h = 1, \dots, H$ **do**
 - 17: Take greedy action $a_h^n = \operatorname{argmax}_a Q_{\hat{w}_h}(x_h^n, a)$ and transit to x_{h+1}^n .
 - 18: **end for**
 - 19: **end for**
-

Definition 2.1 and $\|\phi(x, a)\|_\infty \leq 1$ for any $(x, a) \in \mathcal{X} \times \mathcal{A}$. Assume the learner has oracle access to an exploratory policy π_e defined in Definition 3.1 and $C_{\min}(\Sigma^{\pi_e}, s)$ is a strictly positive universal constant independent of N and d . Choose the regularization parameter $\lambda_1 = H \sqrt{\log(2d)/N}$ and the number of episodes in the exploration phase N_1 as

$$N_1 = \left(\frac{2048s^2H^4N^2}{C_{\min}(\Sigma^{\pi_e}, s)^2} \log(2dH/\delta) \right)^{\frac{1}{3}}.$$

With probability $1 - \delta$, the cumulative regret of online Lasso-FQI satisfies:

$$R_N \leq 2 \left(\frac{2048 \log(2dH/\delta)}{C_{\min}(\Sigma^{\pi_e}, s)^2} \right)^{\frac{1}{3}} H^{\frac{4}{3}} s^{\frac{2}{3}} N^{\frac{2}{3}}. \quad (4.2)$$

With oracle access of an exploratory policy, we obtain a dimension-free sub-linear regret bound. Without oracle access of such an exploratory policy, Theorem 3.3 implies a linear regret lower bound. On the other hand, without considering the sparsity, solving the MDP will suffer linear regret in the high-dimensional regime due to the well-known $\Omega(d\sqrt{N})$ lower bound.

In summary, we emphasize that in high-dimensional regime, exploiting the sparsity to reduce the regret needs an exploratory policy but finding the exploratory policy is as hard as solving the MDP itself - an irresolvable ‘‘chicken and egg’’ problem.

5 COMPARISON WITH CONTEXTUAL BANDITS

In this section we investigate the difference between online RL and linear contextual bandits. When the planning horizon $H = 1$, the episodic MDP becomes to a contextual bandit. Specifically, consider a sparse linear contextual bandit. At n th episode, the environment generates a context x_n i.i.d from a distribution ξ_0 . The learner chooses an action $a_n \in \mathcal{A}$ and receives a reward: $Y_n = \phi(x_n, a_n)^\top \theta + \eta_n$, where $(\eta_n)_{n=1}^N$ is a sequence of independent standard Gaussian random variables and $\theta \in \mathbb{R}^d$ is a s -sparse unknown parameter vector.

We define an analogous exploratory policy as in Definition 3.1: for an exploratory policy π_e in a linear contextual bandit, it will satisfy

$$\sigma_{\min}(\Sigma^{\pi_e}) = \sigma_{\min} \left(\mathbb{E}^{\pi_e} \left[\phi(x_n, a_n) \phi(x_n, a_n)^\top \right] \right) > 0,$$

where $x_n \sim \xi_0$ and $a_n \sim \pi_e(\cdot | x_n)$. In episodic MDPs, since the MDP transition kernel is unknown, we can not find the exploratory policy without solving the MDP. However, in linear contextual bandits, as long as there exists an exploratory policy and the context distribution is known, we can obtain the exploratory policy by solving the following optimization problem:

$$\max_{\pi} \sigma_{\min} \left(\mathbb{E}_{x \sim \xi_0, a \sim \pi(\cdot | x)} \left[\phi(x, a) \phi(x, a)^\top \right] \right).$$

Thus, there is no additional cost of the regret to obtain the exploratory policy. Note that assuming known context distribution is much weaker than assuming known MDP transition kernel since we can learn the context distribution very quickly online. Following the rest step of online Lasso-FQI in Algorithm 1, we can replicate the $\tilde{O}(s^{2/3}N^{2/3})$ regret upper bound without oracle access of the exploratory policy.

6 DISCUSSION

In this paper, we provide the first investigation of online sparse RL in the high-dimensional regime. In general, exploiting the sparsity to minimize the regret is hard without further assumptions. This also highlights some fundamental differences of sparse learning between online RL and supervised learning or contextual bandits.

Acknowledgements

Csaba Szepesvári gratefully acknowledges funding from the Canada CIFAR AI Chairs Program, Amii and NSERC.

References

- Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvári. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pages 1–9, 2012.
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. FLAMBE: Structural complexity and representation learning of low rank MDPs. *arXiv preprint arXiv:2006.10814v2*, 2020.
- András Antos, Csaba Szepesvári, and Rémi Munos. Fitted Q-iteration in continuous action-space MDPs. In *Advances in neural information processing systems*, pages 9–16, 2008.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272, 2017.
- Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294, 2020.
- I. R. Bellman, R. Kalaba, and B. Kotkin. Polynomial approximation – a new computational technique in dynamic programming. *Math. Comp.*, 17(8):155–161, 1963.
- Dimitri P. Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena Scientific, 1995.
- Peter J Bickel, Ya’acov Ritov, Alexandre B Tsybakov, et al. Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*, 2019.
- Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005.
- Matthieu Geist and Bruno Scherrer. ℓ^1 -penalized projected Bellman residual. In *European Workshop on Reinforcement Learning*, pages 89–101. Springer, 2011.
- Matthieu Geist, Bruno Scherrer, Alessandro Lazaric, and Mohammad Ghavamzadeh. A Dantzig selector approach to temporal difference learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 347–354, 2012.
- Mohammad Ghavamzadeh, Alessandro Lazaric, Rémi Munos, and Matthew Hoffman. Finite-sample analysis of Lasso-TD. In *Proceedings of the 28th International Conference on Machine Learning*, pages 1177–1184, 2011.
- Botao Hao, Yaqi Duan, Tor Lattimore, Csaba Szepesvári, and Mengdi Wang. Sparse feature selection makes batch reinforcement learning more sample efficient. *arXiv preprint arXiv:2011.04019*, 2020a.
- Botao Hao, Tor Lattimore, and Mengdi Wang. High-dimensional sparse linear bandits. *Advances in Neural Information Processing Systems*, 33, 2020b.
- Matthew W Hoffman, Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Regularized least squares temporal difference learning with nested ℓ^2 and ℓ^1 penalization. In *European Workshop on Reinforcement Learning*, pages 102–114. Springer, 2011.
- Morteza Ibrahim, Adel Javanmard, and Benjamin V Roy. Efficient reinforcement learning for high dimensional linear quadratic systems. In *Advances in Neural Information Processing Systems*, pages 2636–2644, 2012.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 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*. JMLR.org, 2017.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019.
- Gi-Soo Kim and Myunghee Cho Paik. Doubly-robust lasso bandit. In *Advances in Neural Information Processing Systems*, pages 5869–5879, 2019.
- J Zico Kolter and Andrew Y Ng. Regularization and feature selection in least-squares temporal difference learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 521–528, 2009.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Bo Liu, Sridhar Mahadevan, and Ji Liu. Regularized off-policy TD-learning. In *Advances in Neural Information Processing Systems*, pages 836–844, 2012.
- Christopher Painter-Wakefield and Ronald Parr. Greedy algorithms for sparse reinforcement learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 867–874, 2012.
- Martin L Puterman. *Markov Decision Processes.: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- Zhimei Ren and Zhengyuan Zhou. Dynamic batch learning in high-dimensional sparse linear contextual bandits. *arXiv preprint arXiv:2008.11918*, 2020.
- Paul J Schweitzer and Abraham Seidmann. Generalized polynomial approximations in Markovian decision processes. *Journal of Mathematical Analysis and Applications*, 110(2):568–582, 1985.
- Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, pages 2898–2933, 2019.

- Csaba Szepesvári. *Algorithms for Reinforcement Learning*. Morgan and Claypool, 2010.
- Robert Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Publishing Company, Incorporated, 1st edition, 2008. ISBN 0387790519, 9780387790510.
- Sara A Van De Geer, Peter Bühlmann, et al. On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- Martin J Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- Xue Wang, Mingcheng Wei, and Tao Yao. Minimax concave penalized multi-armed bandit model with high-dimensional covariates. In *International Conference on Machine Learning*, pages 5200–5208, 2018.
- Yining Wang, Yi Chen, Ethan X Fang, Zhaoran Wang, and Runze Li. Nearly dimension-independent sparse linear bandit over small action spaces via best subset selection. *arXiv preprint arXiv:2009.02003*, 2020.
- Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004, 2019.
- Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *International Conference on Machine Learning*, 2020.
- Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirodda, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964, 2020.
- Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted mdps with feature mapping. *arXiv preprint arXiv:2006.13165*, 2020.