

# Computational-Statistical Gap in Reinforcement Learning

**Daniel Kane**

**Sihan Liu**

**Shachar Lovett**

**Gaurav Mahajan**

*University of California, San Diego*

DAKANE@ENG.UCSD.EDU

SIL046@UCSD.EDU

SLOVETT@CS.UCSD.EDU

GMAHAJAN@ENG.UCSD.EDU

**Editors:** Po-Ling Loh and Maxim Raginsky

## Abstract

Reinforcement learning with function approximation has recently achieved tremendous results in applications with large state spaces. This empirical success has motivated a growing body of theoretical work proposing necessary and sufficient conditions under which efficient reinforcement learning is possible. From this line of work, a remarkably simple minimal sufficient condition has emerged for sample efficient reinforcement learning: MDPs with optimal value function  $V^*$  and  $Q^*$  linear in some known low-dimensional features. In this setting, recent works have designed sample efficient algorithms which require a number of samples polynomial in the feature dimension and independent of the size of state space. They however leave finding computationally efficient algorithms as future work and this is considered a major open problem in the community.

In this work, we make progress on this open problem by presenting the first computational lower bound for RL with linear function approximation: unless  $NP=RP$ , no randomized polynomial time algorithm exists for deterministic transition MDPs with a constant number of actions and linear optimal value functions. To prove this, we show a reduction from UNIQUE-SAT, where we convert a CNF formula into an MDP with deterministic transitions, constant number of actions and low dimensional linear optimal value functions. This result also exhibits the first computational-statistical gap in reinforcement learning with linear function approximation, as the underlying statistical problem is information-theoretically solvable with a polynomial number of queries, but no computationally efficient algorithm exists unless  $NP=RP$ . Finally, we also prove a quasi-polynomial time lower bound under the Randomized Exponential Time Hypothesis.

**Keywords:** Machine Learning (cs.LG); Artificial Intelligence (cs.AI); Computational Complexity (cs.CC); Optimization and Control (math.OC); Machine Learning (stat.ML)

## 1. Introduction

Function approximation has a long history in reinforcement learning (Tsitsiklis and Van Roy, 1996; Bertsekas, 2009; Munos and Szepesvári, 2008) and game playing (Shannon, 1950; Tesauro et al., 1995). More recently, this merger of reinforcement learning’s algorithmic techniques with supervised learning’s generalization schemes has achieved tremendous results in various applications with large state spaces, in areas such as game playing (Mnih et al., 2013; Silver et al., 2017; Berner et al., 2019), robotics (Kober et al., 2013) and biology (Senior et al., 2020). Since, one would expect the statistical and computational demand for these algorithms to grow at least linearly with the size of the state space (Jaksch et al., 2010), it is quite surprising that these algorithms generalize so well in large state spaces. That said, the computational requirements for existing algorithms have become exceedingly high. For example, AlphaZero was trained on 5000 tensor processing units (TPUs) for 13 days (Silver et al., 2018) and OpenAI Five trained its DOTA2 bots using 128000

CPUs (Berner et al., 2019) for 180 days (10 months in real time). This leads to a natural fundamental question: are such data and compute requirements fundamental or can we design efficient algorithms for these applications? More generally: *what minimal properties of environments leads to efficient RL algorithms?*

Over the last decade, this question has driven a growing body of theoretical work showing when *sample efficiency* is possible in RL for particular model classes, such as State Aggregation (Li, 2009; Dong et al., 2020), Linear MDPs (Yang and Wang, 2019; Jin et al., 2020), Linear Mixture MDPs (Modi et al., 2020; Ayoub et al., 2020; Zhou et al., 2021), Reactive POMDPs (Krishnamurthy et al., 2016), Block MDPs (Du et al., 2019), FLAMBE (Agarwal et al., 2020), Reactive PSRs (Littman et al., 2001), Linear Bellman Complete (Munos, 2005; Zanette et al., 2020). More generally, there are also a few lines of work which propose general frameworks, consisting of structural conditions which permit sample efficient RL; these include the Bellman rank (Jiang et al., 2016), Witness rank (Sun et al., 2019), Bilinear Classes (Du et al., 2021) and Bellman Eluder (Jin et al., 2021). The goal in these latter works is to develop a unified theory of generalization in RL, analogous to the more classical notions in statistical complexity (e.g. VC-theory and Rademacher complexity) relevant for supervised learning.

A surprisingly minimal assumption which arose from these works is Linear  $Q^*$ & $V^*$  (Du et al., 2021) where both optimal value function  $V^*$  and optimal action-value function  $Q^*$  are linear in some known low-dimensional features. Du et al. (2021) showed that in this setting, there exists sample efficient RL algorithms which *regardless of the number of actions* require a number of samples polynomial in the feature dimension and independent of the size of the state space. However, when only either  $V^*$  or  $Q^*$  are linear, a series of works (Weisz et al., 2020; Wang et al., 2021; Weisz et al., 2021c; Foster et al., 2021) showed that a phase transition occurs as one increases the number of actions: sample efficient algorithms exist for constant number of actions, and quickly transform into information theoretic exponential lower bounds as the number of actions exceeds the dimension of the features underlying  $Q^*$  or  $V^*$ .

Even though we have made considerable progress in understanding the minimal assumptions from the statistical perspective, the computational aspect of this problem is largely unknown. All the settings mentioned above (except under strong assumptions like linear transitions (Jin et al., 2020) and deterministic rewards (Wen and Van Roy, 2013)) do not have computationally efficient algorithms and previous works (Jiang et al., 2016; Du et al., 2021; Weisz et al., 2021a) leave designing computationally efficient algorithms as an important open problem. On the other hand, in spite of failed search for such computationally efficient algorithms over the last few years, there are no computational lower bounds for any of these settings (although previous attempts (Dann et al., 2018) have shown inefficiency of specific algorithms). A case of particular interest is RL under linear function approximation with constant many actions, which includes linear  $Q^*$  (Weisz et al., 2021c), linear  $V^*$  (Weisz et al., 2021d), linear  $Q^*$ & $V^*$  (Du et al., 2021) and linear  $Q^*$ & $V^*$  (reachable states) (Weisz et al., 2021c). In all these settings, we have statistically efficient algorithms when the number of actions are  $O(1)$ , but all the algorithms Du et al. (2020, 2021); Weisz et al. (2021c) take times either exponential in  $d$  or  $H$ . Designing polynomial time algorithms for any of these settings is considered a major open problem in the community.

### 1.1. Our Contributions

In this work, we present the first computational lower bounds for RL with linear function approximation. Before stating our main results, we first need to state some key definitions that we use throughout the paper.

**Markov Decision Process (MDP).** We first define the framework for reinforcement learning, a Markov Decision Process (MDP). We define a deterministic MDP as a tuple  $M = (\mathcal{S}, \mathcal{A}, R, P)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $R : \mathcal{S} \times \mathcal{A} \mapsto \Delta([0, 1])$  is the stochastic reward function<sup>1</sup>, and  $P : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  is the deterministic transition function. An MDP  $M$  defines a discrete time sequential decision process where the agent starts from a starting state  $s_0 \in \mathcal{S}$ . Then, at each time  $t$ , the agent at some current state  $S_t$ , takes action  $A_t$ , receiving reward  $R_t \sim R(S_t, A_t)$  and transitions to next state  $S_{t+1}$ . This goes on till the agent reaches the end state  $\perp$ . Each such trajectory/path from starting state  $s_0$  to end state  $\perp$  is of length at most horizon  $H$ . A deterministic, stationary policy  $\pi : \mathcal{S} \mapsto \mathcal{A}$  specifies a decision-making strategy in which the agent chooses actions adaptively based on the current state, i.e.  $A_t = \pi(S_t)$ . Given a policy  $\pi$  and a state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , the  $Q$ -function and  $V$ -function under a policy  $\pi$  are defined as

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\tau-1} R(S_t, A_t) \mid S_0 = s, \pi \right], \quad Q^\pi(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\tau-1} R(S_t, A_t) \mid S_0 = s, A_0 = a, \pi \right], \quad (1)$$

where  $S_1, A_1, \dots, S_{\tau-1}, A_{\tau-1}$  are obtained by executing policy  $\pi$  in the MDP  $M$  and  $\tau$  is the first time when policy  $\pi$  reaches the end state  $\perp$ , that is  $S_\tau = \perp$  where it always holds that  $\tau \leq H$ . We use  $Q^*$  and  $V^*$  to denote the optimal value functions

$$V^*(s) = \sup_{\pi} V^\pi(s), \quad Q^*(s, a) = \sup_{\pi} Q^\pi(s, a), \quad s \in \mathcal{S}, a \in \mathcal{A}$$

We say that the optimal value functions  $V^*$  and  $Q^*$  can be written as a linear function of  $d$ -dimensional features  $\psi : \mathcal{S} \cup (\mathcal{S} \times \mathcal{A}) \rightarrow \mathbb{R}^d$  if for all state  $s$  and action  $a$ ,  $V^*(s) = \langle \theta, \psi(s) \rangle$  and  $Q^*(s, a) = \langle \theta, \psi(s, a) \rangle$  for some fixed  $\theta \in \mathbb{R}^d$  independent of  $s$  and  $a$ .

**Computational Problems.** We next introduce 3-SAT, a satisfiability problem for 3-CNF formulas. In a 3-SAT problem, we are given as input, a 3-CNF formula  $\varphi$  with  $v$  variables and  $O(v)$  clauses and our goal is to decide if  $\varphi$  is satisfiable. Our computational lower bound is based on a reduction from UNIQUE-3-SAT, a variant of 3-SAT. UNIQUE-3-SAT is the promise version of 3-SAT where the given formula is promised to have either 0 or 1 satisfying assignments.

The focus of this work is the computational RL problem, LINEAR-k-RL. In a LINEAR-k-RL problem with feature dimension  $d$ , we are given access to a deterministic MDP  $M$  with  $k$  actions and horizon  $H = O(d)$  such that the optimal value functions  $Q^*$  and  $V^*$  can be written as a linear function of  $d$ -dimensional features  $\psi$ . Our goal is to output a good policy, which we define as any policy  $\pi$  that satisfies  $V^\pi > V^* - 1/4$ . Note that here  $V^\pi$  and  $V^*$  refers to the value of the policy  $\pi$  and optimal policy respectively at the starting state and is always in  $[0, H]$ <sup>2</sup>. Moreover, the constant  $1/4$  can be replaced by any arbitrary constant  $< 1$ . From now on, we always assume number of actions  $k$  is 2 or 3.

1.  $\Delta([0, 1])$  denotes the set of all distributions over interval  $[0, 1]$ .

2. in our constructions, we satisfy the more stringent condition that  $V^* \in [0, 1]$ .

---

**Complexity problem** LINEAR-k-RL

*Oracle:* a deterministic MDP  $M$  with  $k$  actions, optimal value functions  $V^*$  and  $Q^*$  linear in  $d$  dimensional features  $\psi$  and horizon  $H = O(d)$ .

*Goal:* find policy  $\pi$  such that  $V^\pi > V^* - 1/4$ .

---

We now describe how the algorithm interacts with the MDP. We assume that the algorithm has access to the associated (i) reward function  $R$ , (ii) transition function  $P$  and (iii) features  $\psi$ . For all these functions, the algorithm provides a state  $s$  and action  $a$  (if needed) and receives a random sample from the distribution  $R(s, a)$  (for the reward function), the state  $P(s, a)$  (for the transition function) or feature  $\psi(s)$  or  $\psi(s, a)$  (for the features). We assume that each call accrues constant runtime and input/output for these functions are of size polynomial in feature dimension  $d$ .

We will often talk about randomized algorithm  $A$  solving a problem in time  $t$  with error probability  $p$ . By this we mean (i)  $A$  runs in time  $O(t)$ ; (ii) for satisfiability problems, it returns YES on positive input instances with probability at least  $1 - p$  and returns NO on negative input instances with probability 1; and (iii) for RL problem, it returns a good policy with probability at least  $1 - p$ .

## 1.1.1. NO POLYNOMIAL TIME ALGORITHM FOR LINEAR-2-RL

With these considerations in mind, we present our main result that asserts that unless  $\text{NP}=\text{RP}$ , no randomized polynomial time algorithm can find a good policy in deterministic MDPs with a constant number of actions and linear optimal value functions.

**Theorem 1 (LINEAR-2-RL  $\in$  RP  $\implies$  NP=RP)** *Unless  $\text{NP}=\text{RP}$ , no randomized algorithm can solve LINEAR-2-RL with feature dimension  $d$  in time polynomial in  $d$  with error probability  $1/10$ .*

This resolves the open problem from [Weisz et al. \(2021c\)](#) and [Du et al. \(2021\)](#) by showing that unless  $\text{RP}=\text{NP}$ , no polynomial time randomized algorithm exists for deterministic transition MDPs with a constant number of actions and linear optimal value functions.

Our main technical contribution is a reduction from UNIQUE-3-SAT to LINEAR-3-RL such that a polynomial time algorithm for LINEAR-3-RL implies a polynomial time algorithm for UNIQUE-3-SAT. To achieve this, we use the input for UNIQUE-3-SAT: a 3-CNF formula  $\varphi$  with  $v$  variables, to design an input for LINEAR-3-RL: an MDP  $M_\varphi$  with 3 actions and optimal value functions  $V^*$  and  $Q^*$  linear in  $d$ -dimensional features. On a high level, the MDP is constructed such that each state represents an assignment to the UNIQUE-3-SAT variables and the goal is to “search” for the solution to the UNIQUE-3-SAT instance. In particular, at each state, the 3 actions available to the agent correspond to an unsatisfied clause which ensures at least one action available to the agent decreases the distance to the solution. To incentivize finding the solution, a large reward is awarded on reaching the solution and a very small expected reward on reaching the horizon (this reward is small enough that any polynomial time RL algorithm only receives 0 reward with high probability on reaching the horizon). This ensures that (i) finding a good policy also finds the satisfying assignment of  $\varphi$  and (ii) the optimal value functions  $V^*$  and  $Q^*$  are linear in some low dimensional features. We present this construction in Section 2.

To get lower bounds for LINEAR-2-RL, we use the same construction as above with a small modification. We replace the choice of 3 actions  $a_1, a_2$  and  $a_3$  at every state with a depth-2 binary tree, where the first action is  $a_1$  and the second action leads to a new state which has actions  $a_2$  and

$a_3$ . This allows us to simulate the hard 3-action MDP using a 2-action MDP while increasing our feature dimension  $d$  by at most a quadratic factor. We present this construction in Section 3.

These reductions allow us to simulate a polynomial time algorithm for UNIQUE-3-SAT on input  $\varphi$  by running the polynomial time algorithm for LINEAR-2-RL on MDP  $M_\varphi$ . More formally, our reduction gives a polynomial relationship between the complexity of UNIQUE-3-SAT and LINEAR-2-RL: a polynomial  $d^q$  time algorithm for LINEAR-2-RL implies a polynomial  $v^{O(q^2)}$  time algorithm for UNIQUE-3-SAT.

**Proposition 2** *Suppose  $q \geq 1$ . If LINEAR-2-RL with feature dimension  $d$  can be solved in time  $d^q$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $v^{O(q^2)}$  with error probability  $1/8$ .*

This relates the complexity of UNIQUE-3-SAT to LINEAR-2-RL and LINEAR-3-RL. To relate these problems to complexity class NP, we use a seminal result from Valiant and Vazirani (1985) which showed that uniqueness of solution can not be used to solve search problems quickly. In particular, they showed a randomized polynomial time reduction from 3-SAT to UNIQUE-3-SAT.

**Theorem 3 (Valiant-Vazirani Theorem)** *Unless  $\text{NP}=\text{RP}$ , no polynomial time randomized algorithm can solve UNIQUE-3-SAT with error probability  $1/8$ .*

Combining our reduction with Valiant-Vazirani Theorem proves our main result—Theorem 1.

### 1.1.2. QUASI-POLYNOMIAL LOWER BOUND FOR LINEAR-2-RL

We now present computational lower bound under a strengthening of  $\text{NP} \neq \text{RP}$  conjecture, Randomized Exponential Time Hypothesis (rETH) (Dell et al., 2014), which asserts that probabilistic algorithms can not decide if a given 3-SAT problem with  $v$  variables and  $O(v)$  clauses is satisfiable in sub-exponential time.

**Definition 4 (Randomized Exponential Time Hypothesis (rETH))** *There is a constant  $c > 0$  such that no randomized algorithm can decide 3-SAT with  $v$  variables in time  $2^{cv}$  with error probability  $1/2$ .*

Randomized Exponential Time Hypothesis along with many variants motivated by Exponential Time Hypothesis (Impagliazzo and Paturi, 2001) have been influential in discovering hardness results for a variety of problems see, e.g. Cygan et al. (2015); Williams (2019). Under Randomized Exponential Time Hypothesis, our main result is a quasi-polynomial computational lower bound for learning good policies in deterministic MDPs with linear optimal value functions.

**Theorem 5 (Quasi-polynomial lower bound for LINEAR-2-RL)** *Under rETH, no randomized algorithm can solve LINEAR-2-RL with feature dimension  $d$  in time  $d^{O(\log d / \log \log d)}$  with error probability  $1/10$ .*

This improves over our super-polynomial lower bound albeit depending on a much stronger hardness assumption. To prove this result, we use a different choice of parameters in our reduction and set the feature dimension  $d$  to be sub-exponential in the number of variables  $v$  to get the following:

**Proposition 6** *If LINEAR-2-RL with feature dimension  $d$  can be solved in time  $d^{O(\log d / \log \log d)}$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $2^{O(v / \log v)}$  with error probability  $1/8$ .*

Here its important to note that we can not use Valiant-Vazirani Theorem to relate UNIQUE-3-SAT and 3-SAT, since it is consistent with Valiant-Vazirani Theorem that UNIQUE-3-SAT is solvable in  $2^{\sqrt{v}}$  time but 3-SAT takes  $2^v$  time. Therefore, we use a more refined lower bound for UNIQUE-3-SAT from Calabro et al. (2008) which showed that if UNIQUE-3-SAT with  $v$  variables can be solved in time  $2^{\alpha v}$  for every  $\alpha > 0$ , then so can  $k$ -SAT for all  $k \geq 3$ .

**Theorem 7 (Calabro et al. (2008))** *Assuming rETH is true, there exists a constant  $c > 0$  such that no randomized algorithm can solve UNIQUE-3-SAT with  $v$  variables in time  $2^{cv}$  with error probability  $1/2$ .*

In conjunction with our reduction, this gives a quasi-polynomial lower bound for LINEAR-2-RL under rETH. We leave as an open problem if the techniques introduced in this work can be used to prove an exponential lower bound for LINEAR-2-RL under rETH.

Our results give evidence that even though having linear optimal value functions is sufficient for sample efficient reinforcement learning Du et al. (2021), it is not sufficient for computationally efficient reinforcement learning. More assumptions are required for computationally efficient algorithms, in addition to optimal value functions  $Q^*$  and  $V^*$  being linear in low-dimensional features, for example sub-optimality gap (Du et al., 2020). We hope that this work will open up new research avenues for finding minimal sufficient conditions for computationally efficient reinforcement learning. We now discuss a few further notable implications of this work.

- *Computational-Statistical Gap:* There are many problems which exhibit computational-statistical gaps i.e. regimes where the underlying statistical problem is information theoretically possible but no computationally efficient algorithm exists. Examples include community detection (Holland et al., 1983; McSherry, 2001; Abbe and Sandon, 2015), planted clique (Alon et al., 1998; Barak et al., 2019) and sparse principal component analysis (Berthet and Rigollet, 2013a,b). To the best of our knowledge, our computational lower bound is the first computational-statistical gap in reinforcement learning with function approximation. When both optimal value functions  $Q^*$  and  $V^*$  are linear, MDPs with any number of actions are statistically easy to solve (Du et al., 2021) but our results show that no polynomial time algorithm can solve these MDPs even with a constant number of actions, unless NP=RP.
- *Natural Problem in NP \ P:* There has been quite a lot of recent work in complexity theory literature on proving quasi-polynomial lower bounds based on Exponential Time Hypothesis (for e.g. dense constraint satisfaction problems (Aaronson et al., 2014), approximating best nash equilibrium (Braverman et al., 2015) and approximating densest  $k$ -subgraph with perfect completeness (Braverman et al., 2017)). This work adds RL with deterministic transition, linear bounded optimal value functions  $V^*$ ,  $Q^*$  and constant number of actions as another natural problem in NP but not in P unless NP=RP.

**Remainder of this paper.** In Section 2 and Section 3, we present our lower bound constructions for 3 action and 2 action MDPs respectively.

## 2. Lower Bound for MDPs with 3 actions

In this section, we will prove the reduction, Theorem 8 and Theorem 9, restated versions of Theorem 2 and Theorem 6 for LINEAR-3-RL. The overall idea is to first build a randomized algorithm  $\mathcal{A}_{SAT}$  which can decide UNIQUE-3-SAT using a randomized algorithm  $\mathcal{A}_{RL}$  which solves LINEAR-3-RL. The two reductions only differ in their settings of parameters.

In the first setting, which we use to prove that no polynomial time algorithm exists for LINEAR-3-RL, we set the feature dimension  $d$  to be polynomial in the number of variables  $v$ . Under this setting, we can build a polynomial time randomized algorithm for UNIQUE-3-SAT using a polynomial time randomized algorithm for LINEAR-3-RL.

**Proposition 8** *Suppose  $q \geq 1$ . If LINEAR-3-RL with feature dimension  $d$  can be solved in time  $d^q$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $O(v^{8q+16q^2})$  with error probability  $1/8$ .*

In the second setting, which we use to prove a quasi-polynomial lower bound for LINEAR-3-RL, we set the feature dimension  $d$  to be sub-exponential in the number of variables  $v$ . This allows us to transform an exponential time lower bound for UNIQUE-3-SAT into a quasi-polynomial lower bound for LINEAR-3-RL.

**Proposition 9** *If LINEAR-3-RL with feature dimension  $d$  can be solved in time  $d^{\log d / (32 \log \log d)}$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $2^{O(v/\log v)}$  with error probability  $1/8$ .*

Before we prove these results, we give a brief outline of our reduction from UNIQUE-3-SAT to LINEAR-3-RL. On a high level, we construct an MDP where the goal is to "search" for the solution  $w^*$  to a UNIQUE-3-SAT instance with  $v$  variables. In particular, at each time, the agent is given an unsatisfied clause and asked to flip assignment for a variable present in the clause. Notice that since the clause is unsatisfied, there must be at least one variable whose assignment differs from the solution and therefore, the agent can "reach" the solution in at most  $d(w, w^*)$  steps. To incentivize the agent, if the agents at time  $l$  finds the solution i.e.  $w = w^*$  or reaches the end of the MDP i.e.  $l = H$ , it receives reward according to the following degree- $r$  polynomial

$$g(l, w) = \left(1 - \frac{l + \text{dist}(w, w^*)}{H + v}\right)^r.$$

We show how to build an MDPs from a UNIQUE-3-SAT instance in Section 2.1. Furthermore, we show that the optimal value functions  $V^*$  and  $Q^*$  for the constructed MDP are linear in  $d = O(v^r)$ -dimensional features. Since the expected reward at last layer of the MDP is  $O(v^{-r^2})$  (which can be replaced with 0 for any poly( $d$ ) time RL algorithm), the only non-zero reward is achieved by solving the underlying UNIQUE-3-SAT instance, proving our reduction. We give a formal argument in Section 2.2, where we show how to build a randomized algorithm for UNIQUE-3-SAT using a randomized algorithm for LINEAR-3-RL. In Section 2.3, we discuss the two different settings of parameters which will prove Theorem 8 and Theorem 9.

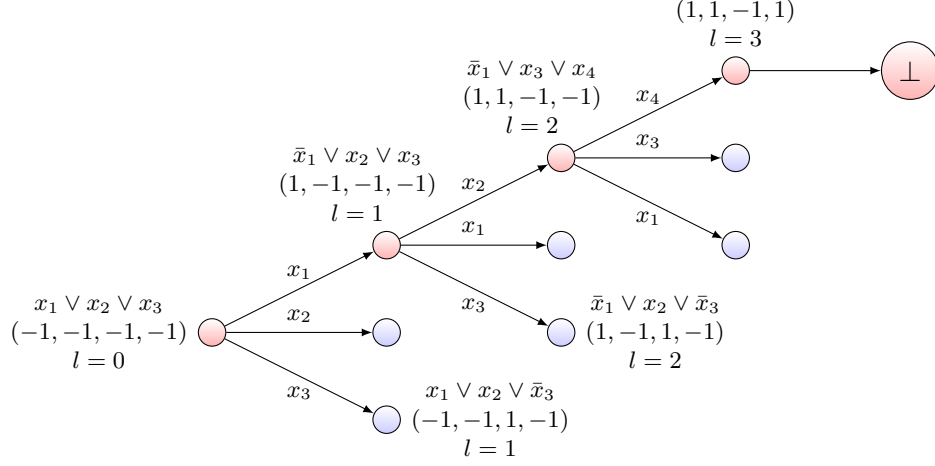


Figure 1: Example construction of 3-action MDP  $M_\varphi$  from a 3-CNF formula  $(x_1 \vee x_2 \vee x_3) \wedge (\bar{x}_1 \vee x_2 \vee x_3) \wedge (\bar{x}_1 \vee x_3 \vee x_4) \wedge (x_1 \vee x_2 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee x_2 \vee \bar{x}_3) \wedge (\bar{x}_3 \vee \bar{x}_3 \vee \bar{x}_3) \wedge (x_1 \vee x_1 \vee x_1)$ . The only satisfying assignment for this formula is  $(1, 1, -1, 1)$ . The states are labelled by the corresponding assignment and unsatisfied clause which decides the available actions. The states in the optimal path are colored in red.

### 2.1. From 3-CNF formulas to 3-action MDPs

We will start by defining a mapping from an input of UNIQUE-3-SAT problem: 3-CNF formula  $\varphi$  with  $v$  variables and  $O(v)$  clauses to an MDP  $M_\varphi$  with 3 actions and  $H = O(d)$  horizon with optimal value functions linear in  $d$  dimensions. Our informal goal is to design an MDP  $M_\varphi$  such that finding a good policy also implies finding the satisfying assignment for the formula  $\varphi$ . We now formally describe the MDP  $M_\varphi$  when the formula  $\varphi$  has a unique satisfying assignment  $w^* \in \{-1, 1\}^v$  and later show how the MDP  $M_\varphi$  differs when the formula  $\varphi$  has no solution. See Figure 1 for an example.

**Transitions.** In our setting, it will be useful to visualize an MDP as a tree, where nodes represent states and edges represent actions. A policy is then a sequence of actions or equivalently a path in the aforementioned tree. The MDP  $M_\varphi$  is a ternary tree i.e. each state/node in the tree has 3 children. The transitions/dynamics are deterministic i.e. the first action goes to first child, the second action goes to second child and so on.

**Assignments.** Each state is associated with an assignment to the  $v$  variables i.e. a binary vector in  $\{-1, 1\}^v$  and a natural number  $l$  denoting the depth of the state. Our goal here is to choose assignments such that it is always possible to choose an action which decreases the hamming distance to the satisfying assignment. The root in the tree is associated with the all zeroes assignment  $(-1, -1, \dots, -1)$ . For any state  $s$  with a non-satisfying assignment  $w = (w_1, w_2, \dots, w_v) \neq w^*$ , the assignment associated to the three children are as follows. Since  $w$  is not a satisfying assignment, consider the first unsatisfied clause with variables  $x_{i_1}, x_{i_2}, x_{i_3}$ . The first child is associated with the assignment where the  $i_1$ -th bit of  $w$  is flipped, the second child is associated with vector where  $i_2$ -th bit is flipped and so on. More formally, the assignment associated to  $j$ -th child is



$(w'_1, w'_2, \dots, w'_v)$  where  $w'_k = \neg w_k$  if  $k = i_j$  and  $w'_k = w_k$  otherwise. The two exceptions to this are (i) states with the satisfying assignment  $w^*$  and (ii) states at the last level  $H$ . For such states, all actions go to the end state  $\perp$ .

**Rewards.** To ensure that finding good policies implies finding the satisfying assignment in our MDP, we will only give rewards when a satisfying assignment is found or at the last layer. More formally, the rewards everywhere are zero except on (i) states with the satisfying assignment  $w^*$  and (ii) states on the last level  $H$ . In both the cases above, say the state is at level  $l$  with assignment  $w$ , then the associated reward distribution for any action is a Bernoulli distribution  $Ber(g(l, w))$  where

$$g(l, w) = \left(1 - \frac{l + \text{dist}(w, w^*)}{H + v}\right)^r$$

and the Bernoulli distribution  $Ber(p)$  is 1 with probability  $p$  and 0 with probability  $1 - p$ . Here  $r$  is a parameter which we will specify in Section 2.3. When the formula  $\varphi$  has no satisfying assignment, all rewards are 0. Note that in our simulation (Section 2.2), we don't know/use  $w^*$  and instead use an approximate reward function that is easy to compute.

**Linear Optimal Value Functions.** We next show that in the MDP  $M_\varphi$ , the optimal value functions  $V^*$  and  $Q^*$  can be written as a linear function of  $d = O(v^r)$  dimensional features  $\psi$ , where  $\psi(s)$  or  $\psi(s, a)$  depends only on  $w$ , the corresponding assignment, and  $l$ , the depth of the state.

**Proposition 10** *For any state  $s$  in level  $l$  with assignment  $w$  and action  $a$ ,*

- (i) *the optimal value function is  $V^*(s) = g(l, w)$ .*
- (ii) *for large enough  $v$ , there exists features  $\psi(s), \psi(s, a) \in \mathbb{R}^d$  with feature dimension  $d \leq 2v^r$  depending only on state  $s$  and action  $a$ ; and  $\theta \in \mathbb{R}^d$  depending only on  $w^*$  such that  $V^*$  and  $Q^*$  can be written as a linear function of features  $\psi$  i.e.  $V^*(s) = \langle \theta, \psi(s) \rangle$  and  $Q^*(s, a) = \langle \theta, \psi(s, a) \rangle$ .*

We present the proof in Appendix A. Even though  $\psi(s)$  does not depend on  $w^*$ , unlike the constructions of Weisz et al. (2021b,c); Wang et al. (2021),  $\psi(s)$  does depend on the MDP  $M_\varphi$  making this construction statistically easy but computationally hard to solve.

## 2.2. From RL algorithms to 3-SAT algorithms

We now build a randomized algorithm  $\mathcal{A}_{SAT}$  for UNIQUE-3-SAT using a randomized algorithm  $\mathcal{A}_{RL}$  for the RL problem. However, as mentioned before, since the runtime for  $\mathcal{A}_{RL}$  accrues only constant runtime for each call to the MDP oracle, to efficiently build  $\mathcal{A}_{SAT}$  using  $\mathcal{A}_{RL}$ , we need to be able to efficiently simulate the calls to MDP oracle, namely: calls to the reward function, the transition function and the features. To do so, we build an ‘‘approximate’’ simulator  $\bar{M}_\varphi$  for the MDP oracle  $M_\varphi$ . The simulator  $\bar{M}_\varphi$  is exactly MDP  $M_\varphi$  in terms of transition function and features associated with the MDP  $M_\varphi$ , but differs in the reward function at the last layer which is always 0 for the simulator  $\bar{M}_\varphi$ . This modification is crucial for an efficient reduction because unlike transitions and features for any state which can be computed in time  $\text{poly}(d)$  on the MDP  $M_\varphi$ , the rewards at the last layer when  $\text{dist}(w, w^*) \neq 0$  require access to  $w^*$  which can not be done efficiently. With the purposed modification, we can execute each call to simulator  $\bar{M}_\varphi$  in time  $\text{poly}(d)$ .

**Algorithm.** On input 3-CNF formula  $\varphi$ ,  $\mathcal{A}_{SAT}$  runs the algorithm  $\mathcal{A}_{RL}$  replacing each call to MDP oracle  $M_\varphi$  with the corresponding call to simulator  $\bar{M}_\varphi$ . Recall that the output for the RL algorithm in our setting is a sequence of actions. If the sequence of actions returned by  $\mathcal{A}_{RL}$  ends on a state with assignment  $w$ ,  $\mathcal{A}_{SAT}$  outputs YES if  $w$  is the satisfying assignment and returns NO otherwise.

**Correctness.** We set the horizon  $H = v^r$ . We will assume throughout that  $r \geq 2$  and that the runtime of  $\mathcal{A}_{RL}$  is  $\leq v^{r^2/4}$ . Different settings of  $r$  satisfying these assumptions will prove Theorem 8 and Theorem 9 for 3-action MDPs, which we will discuss in Section 2.3. To complete our reduction, we will show the following:

- (i) If algorithm  $\mathcal{A}_{RL}$  outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$ , then  $\mathcal{A}_{SAT}$  on 3-CNF formula  $\varphi$  outputs YES if  $\varphi$  is satisfiable and NO otherwise.
- (ii) If  $\mathcal{A}_{RL}$  with access to MDP oracle  $M_\varphi$  outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$  with error probability  $1/10$ , then  $\mathcal{A}_{RL}$  with access to simulator  $\bar{M}_\varphi$  outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$  with error probability  $1/8$ .

These together will show that  $\mathcal{A}_{SAT}$  solves UNIQUE-3-SAT with error probability  $\leq 1/8$ . We start by proving that if  $\mathcal{A}_{RL}$  succeeds on MDP  $\bar{M}_\varphi$ , then  $\mathcal{A}_{SAT}$  succeeds on 3-CNF formula  $\varphi$ . This follows from the fact that any good policy in the MDP  $M_\varphi$  must reach a state with satisfying assignment  $w^*$ .

**Proposition 11** *Suppose  $r > 1$  and horizon  $H = v^r$ . If  $\mathcal{A}_{RL}$  outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$ , then  $\mathcal{A}_{SAT}$  on 3-CNF formula  $\varphi$  outputs YES if  $\varphi$  is satisfiable and NO otherwise.*

**Proof** Since algorithm  $\mathcal{A}_{SAT}$  always returns NO on an unsatisfiable formula, we restrict our attention to a satisfiable formula  $\varphi$ . In the MDP  $M_\varphi$ , (i) rewards are “very small” everywhere except on reaching the satisfying assignment i.e. the expected reward at the last layer in the MDP  $M_\varphi$  is upper bounded by (for large enough  $v$  and  $r > 1$ )

$$\left(1 - \frac{H}{H+v}\right)^r = \left(\frac{v}{H+v}\right)^r \leq v^{-r^2+r} < 1/4$$

and (ii) the optimal value  $V^*$  is large

$$V^* \geq \left(1 - \frac{v}{H+v}\right)^r = \left(1 + \frac{v}{v^r}\right)^{-r} \geq 1 - \frac{rv}{v^r} \geq \frac{1}{2}$$

where the second last inequality follows from Bernoulli’s inequality and the last inequality holds for large enough  $v$  and  $r > 1$ . Therefore, if the value of policy is large i.e.  $V^\pi > V^* - 1/4$ , then the policy  $\pi$  (and therefore the corresponding sequence of actions) has to end on a state with the satisfying assignment  $w^*$ . By construction of  $\mathcal{A}_{SAT}$ , this implies  $\mathcal{A}_{SAT}$  will succeed on the formula  $\varphi$ . ■

Since we can not simulate the rewards on MDP oracle  $M_\varphi$  efficiently, our reduction runs the algorithm  $\mathcal{A}_{RL}$  on an approximate simulator  $\bar{M}_\varphi$ . However, it’s not clear why  $\mathcal{A}_{RL}$  would still succeed when each call to MDP oracle is replaced by a call to the simulator  $\bar{M}_\varphi$ . The following proposition shows that in fact  $\mathcal{A}_{RL}$  would succeed on the outputs of simulator  $\bar{M}_\varphi$  albeit with a smaller constant probability.

**Proposition 12** *Suppose  $r \geq 2$  and horizon  $H = v^r$ . Suppose  $\mathcal{A}_{RL}$  with access to MDP oracle  $M_\varphi$  runs in time  $v^{r^2/4}$  and outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$  with error probability  $1/10$ . Then  $\mathcal{A}_{RL}$  with access to simulator  $\bar{M}_\varphi$ , still running in time  $v^{r^2/4}$ , outputs a policy  $\pi$  such that  $V^\pi > V^* - 1/4$  with error probability  $1/8$ .*

**Proof** Let  $\Pr_{M_\varphi}$  and  $\Pr_{\bar{M}_\varphi}$  denote the distribution on the observed rewards and output policies induced by the algorithm  $\mathcal{A}_{RL}$  when running on access to MDP oracle  $M_\varphi$  and simulator  $\bar{M}_\varphi$  respectively. Let  $R_i$  denote the reward received on the last layer at the end of  $i$ -th trajectory. Let  $T$  be the total number of trajectories sampled by algorithm  $\mathcal{A}_{RL}$  when running on access to MDP oracle  $M_\varphi$ . By our assumption,  $\mathcal{A}_{RL}$  runs in time  $v^{r^2/4}$  and therefore  $T \leq v^{r^2/4}$ . Since the expected reward at the last layer in the MDP  $M_\varphi$  is upper bounded by (for large enough  $v$  and  $r \geq 2$ )

$$\left(1 - \frac{H}{H+v}\right)^r = \left(\frac{v}{H+v}\right)^r \leq v^{-r^2+r} \leq v^{-\frac{r^2}{2}}$$

and the algorithm only visits at most  $v^{r^2/4}$  states on last layer, we get by the union bound that with high probability all the rewards at the last level are zero. More precisely (and assuming  $v$  is large enough),

$$\Pr_{M_\varphi} [R_i = 0 \forall i \in [T]] \geq 1 - v^{-r^2/4} \geq \frac{4}{5}$$

We say  $\mathcal{A}_{RL}$  succeeds with access to  $M_\varphi$  (or  $\bar{M}_\varphi$ ) if the output policy  $\pi$  after running for time at most  $v^{r^2/4}$  satisfies  $V^\pi > V^* - 1/4$ . Using the above reasoning and the assumption that  $\mathcal{A}_{RL}$  succeeds with access to MDP oracle  $M_\varphi$  with probability  $9/10$  implies

$$\Pr_{M_\varphi} [\mathcal{A}_{RL} \text{ succeeds with access to } M_\varphi \mid R_i = 0 \forall i \in [T]] \geq \frac{\frac{9}{10} - \frac{1}{5}}{\frac{4}{5}} = \frac{7}{8}$$

Note that the marginal distributions  $\Pr_{M_\varphi}$  and  $\Pr_{\bar{M}_\varphi}$  on output policy  $\pi$  given  $R_i = 0 \forall i \in [T]$  are exactly the same because MDP oracle  $\bar{M}_\varphi$  and simulator  $M_\varphi$  only differ on last layer rewards. This implies

$$\begin{aligned} & \Pr_{\bar{M}_\varphi} [\mathcal{A}_{RL} \text{ succeeds with access to } \bar{M}_\varphi \mid R_i = 0 \forall i \in [T]] \\ &= \Pr_{M_\varphi} [\mathcal{A}_{RL} \text{ succeeds with access to } M_\varphi \mid R_i = 0 \forall i \in [T]] \end{aligned}$$

Since,  $\Pr_{\bar{M}_\varphi} [R_i = 0 \forall i \in [T]] = 1$ , we conclude that

$$\Pr_{\bar{M}_\varphi} [\mathcal{A}_{RL} \text{ succeeds with access to } \bar{M}_\varphi] \geq \frac{7}{8}$$

■

### 2.3. Setting of Parameters

It follows from Theorems 10 to 12 that if LINEAR-3-RL with feature dimension  $d = 2v^r$  can be solved in time  $v^{r^2/4}$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $d \cdot v^{r^2/4}$  with error probability  $1/8$  (here the extra  $d$  factor is because each call to the simulator  $\bar{M}_\varphi$  takes  $d$  time). As we increase  $r$ , we decrease the expected reward available to the algorithm at the last layer on the order of  $v^{-O(r^2)}$ , making the problem harder. However, increasing  $r$  also increases the feature dimension on the order of  $v^r$ . This non-polynomial gap in the feature dimension and expected reward at the last layer will give our main reduction.

We use two different settings of  $r$  to prove our lower bounds: (i)  $r$  is constant wrt number of variables  $v$  and (ii)  $r^2$  is almost linear in the number of variables  $v$ . In Appendix A, we show how the first setting proves Theorem 8 and the second setting proves Theorem 9. We also prove a more general version, Theorem 14, which shows that a quasi-polynomial algorithm for LINEAR-3-RL implies a quasi-polynomial algorithm for UNIQUE-3-SAT.

## 3. Lower Bound for MDPs with 2 actions

In this section, we prove computational lower bound for LINEAR-2-RL. Similar to Section 2, our proof is based on reduction from UNIQUE-3-SAT. We will modify the MDP  $M_\varphi$  with three actions into  $\tilde{M}_\varphi$  by introducing some intermediate states. See Figure 3 for an example of this modification for a single state.

**Intermediate states.** Recall that in  $M_\varphi$  each state is associated with an assignment. Let the  $i$ -th clause, which consists of three variables  $x_{i_1}, x_{i_2}, x_{i_3}$ , be the first unsatisfied clause. Then, the three actions available each correspond to flipping one of the variable in the clause. We will replace them by two actions: while one action still flips the last variable  $x_{i_3}$ , the other action leads to an intermediate state  $s_{[i_1, i_2]}$ . At the state  $s_{[i_1, i_2]}$ , two actions are available: one flips  $x_{i_1}$  and the other flips  $x_{i_2}$ .

**Depth of state.** In the 3 action MDP  $M_\varphi$ , the depth of a state is simply the length of the path that ends at the state. Here, we define the depth to be the number of non-intermediate states included in the path. That being said, the intermediate states will have the same depth as their parents.

**Rewards.** The rewards are the same as those in the 3 action MDP. Namely, rewards are only given at last layer or when the assignment is  $w^*$ . In particular, for a state with assignment  $w$  and depth  $l$ , the reward distribution is  $Ber(g(l, w))$  where

$$g(l, w) = \left(1 - \frac{l + \text{dist}(w, w^*)}{H + v}\right)^r.$$

We now show that, even with this modification, the optimal value functions  $V^*$  and  $Q^*$  can still be written as a linear function of some low dimensional features.

**Proposition 13** *For any state  $s$  in level  $l$  with assignment  $w$  and action  $a$ ,*

- (i) *If  $s$  is a non-intermediate state, then the optimal value function is  $V^*(s) = g(l, w)$ .*

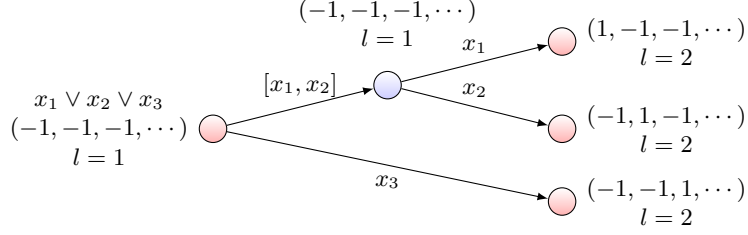


Figure 2: Part of a 2-action MDP corresponding to the CNF clause  $(x_1 \vee x_2 \vee x_3)$ . The non-intermediate states are colored red and the intermediate states are colored blue.

(ii) If  $s$  is an intermediate state that leads to actions which flip coordinates  $i_1$  and  $i_2$ , then the optimal value function is

$$V^*(s_{[i_1, i_2]}) = \left( 1 - \frac{l + \text{dist}(w, w^*) + 2 \cdot \mathbb{1}\{w_{i_1} = w_{i_1}^*\} \cdot \mathbb{1}\{w_{i_2} = w_{i_2}^*\}}{H + v} \right)^r.$$

(iii) for feature dimension  $d = 2v^{2r}$ , there exists features  $\psi(s), \psi(s, a) \in \mathbb{R}^d$  depending only on state  $s$  and action  $a$ ; and  $\theta \in \mathbb{R}^d$  depending only on  $w^*$  such that  $V^*$  and  $Q^*$  can be written as a linear function of features  $\psi$  i.e.  $V^*(s) = \langle \theta, \psi(s) \rangle$  and  $Q^*(s, a) = \langle \theta, \psi(s, a) \rangle$ .

By Theorem 13, the feature dimension  $d$  of the 2 action MDP and the number of variables  $v$  in the UNIQUE-3-SAT instance are related by  $d = 2v^{2r} \leq v^{3r}$ . We present proof of Theorem 13 and how setting of  $r$  proves Theorem 2 and Theorem 6 in Appendix B.

## References

- Scott Aaronson, Russell Impagliazzo, and Dana Moshkovitz. Am with multiple merlins. *2014 IEEE 29th Conference on Computational Complexity (CCC)*, pages 44–55, 2014.
- Emmanuel Abbe and Colin Sandon. Detection in the stochastic block model with multiple clusters: proof of the achievability conjectures, acyclic bp, and the information-computation gap. *arXiv preprint arXiv:1512.09080*, 2015.
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *arXiv preprint arXiv:2006.10814*, 2020.
- Noga Alon, Michael Krivelevich, and Benny Sudakov. Finding a large hidden clique in a random graph. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, page 594–598, 1998.
- Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin F Yang. Model-based reinforcement learning with value-targeted regression. *arXiv:2006.01107*, 2020.
- Boaz Barak, Samuel B. Hopkins, Jonathan A. Kelner, Pravesh Kothari, Ankur Moitra, and Aaron Potechin. A nearly tight sum-of-squares lower bound for the planted clique problem. *SIAM J. Comput.*, 48:687–735, 2019.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. *CoRR*, 2019.
- Quentin Berthet and Philippe Rigollet. Complexity theoretic lower bounds for sparse principal component detection. In Shai Shalev-Shwartz and Ingo Steinwart, editors, *Proceedings of the 26th Annual Conference on Learning Theory*, volume 30, pages 1046–1066, 2013a.
- Quentin Berthet and Philippe Rigollet. Optimal detection of sparse principal components in high dimension. *The Annals of Statistics*, 41(4):1780–1815, 2013b.
- Dimitri P. Bertsekas. *Neuro-Dynamic Programming*, pages 2555–2560. Springer US, Boston, MA, 2009.
- Mark Braverman, Young Kun Ko, and Omri Weinstein. Approximating the best nash equilibrium in  $n^{o(\log n)}$ -time breaks the exponential time hypothesis. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, page 970–982, 2015.
- Mark Braverman, Young Kun-Ko, Aviad Rubinfeld, and Omri Weinstein. Eth hardness for densest-k-subgraph with perfect completeness. *SODA '17*, page 1326–1341, USA, 2017. Society for Industrial and Applied Mathematics.
- Chris Calabro, Russell Impagliazzo, Valentine Kabanets, and Ramamohan Paturi. The complexity of unique k-sat: An isolation lemma for k-cnfs. *Journal of Computer and System Sciences*, 74(3):386–393, 2008. Computational Complexity 2003.

- Marek Cygan, Fedor Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michał Pilipczuk, and Saket Saurabh. *Lower Bounds Based on the Exponential-Time Hypothesis*, pages 467–521. 07 2015. ISBN 978-3-319-21274-6. doi: 10.1007/978-3-319-21275-3\_14.
- Christoph Dann, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. On oracle-efficient pac rl with rich observations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, page 1429–1439, 2018.
- Holger Dell, Thore Husfeldt, Dániel Marx, Nina Taslamán, and Martin Wahlén. Exponential time complexity of the permanent and the tutte polynomial. *ACM Trans. Algorithms*, 10(4), 2014.
- Shi Dong, Benjamin Van Roy, and Zhengyuan Zhou. Provably efficient reinforcement learning with aggregated states, 2020.
- Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in rl. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 2826–2836, 18–24 Jul 2021.
- Simon S Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudík, and John Langford. Provably efficient RL with rich observations via latent state decoding. In *International Conference on Machine Learning*, 2019.
- Simon S Du, Jason D Lee, Gaurav Mahajan, and Ruosong Wang. Agnostic Q-learning with function approximation in deterministic systems: Tight bounds on approximation error and sample complexity. In *Advances in Neural Information Processing Systems*, 2020.
- Dylan J. Foster, Sham M. Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making, 2021.
- Paul W. Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5(2):109–137, 1983.
- Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *J. Comput. Syst. Sci.*, 62(2):367–375, 2001.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. Contextual decision processes with low bellman rank are pac-learnable, 2016.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, 2020.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *arXiv preprint arXiv:2102.00815*, 2021.
- 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.

- Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Pac reinforcement learning with rich observations. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 1848–1856, 2016.
- Lihong Li. *A Unifying Framework for Computational Reinforcement Learning Theory*. PhD thesis, USA, 2009. AAI3386797.
- Michael L Littman, Richard S Sutton, and Satinder P Singh. Predictive representations of state. In *NIPS*, volume 14, page 30, 2001.
- F. McSherry. Spectral partitioning of random graphs. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pages 529–537, 2001.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *Conference on Artificial Intelligence and Statistics*, 2020.
- Rémi Munos. Error bounds for approximate value iteration. In *Proceedings of the National Conference on Artificial Intelligence*, volume 20, page 1006. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.
- Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9(5), 2008.
- Andrew W. Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Žídek, Alexander W. R. Nelson, Alex Bridgland, Hugo Penedones, Stig Petersen, Karen Simonyan, Steve Crossan, Pushmeet Kohli, David T. Jones, David Silver, Koray Kavukcuoglu, and Demis Hassabis. Improved protein structure prediction using potentials from deep learning. *Nature*, 2020.
- Claude E Shannon. Xxii. programming a computer for playing chess. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 41(314):256–275, 1950.
- 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.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 2018.
- 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*, 2019.



- Gerald Tesauro et al. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- John N Tsitsiklis and Benjamin Van Roy. Feature-based methods for large scale dynamic programming. *Machine Learning*, 22(1):59–94, 1996.
- L G Valiant and V V Vazirani. Np is as easy as detecting unique solutions. In *Proceedings of the Seventeenth Annual ACM Symposium on Theory of Computing*, page 458–463, 1985.
- Yuanhao Wang, Ruosong Wang, and Sham M. Kakade. An exponential lower bound for linearly-realizable mdps with constant suboptimality gap, 2021.
- Gellert Weisz, Philip Amortila, and Csaba Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions, 2020.
- Gellért Weisz, Philip Amortila, Barnabás Janzer, Yasin Abbasi-Yadkori, Nan Jiang, and Csaba Szepesvári. On query-efficient planning in mdps under linear realizability of the optimal state-value function, 2021a.
- Gellért Weisz, Philip Amortila, and Csaba Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. In *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, volume 132, pages 1237–1264. PMLR, 2021b.
- Gellért Weisz, Csaba Szepesvári, and András György. Tensorplan and the few actions lower bound for planning in mdps under linear realizability of optimal value functions. *arXiv preprint arXiv:2110.02195*, 2021c.
- Gellért Weisz, Philip Amortila, Barnabás Janzer, Yasin Abbasi-Yadkori, Nan Jiang, and Csaba Szepesvári. On query-efficient planning in mdps under linear realizability of the optimal state-value function, 2021d.
- Zheng Wen and Benjamin Van Roy. Efficient exploration and value function generalization in deterministic systems. In *Advances in Neural Information Processing Systems*, 2013.
- Virginia Vassilevska Williams. On some fine-grained questions in algorithms and complexity. *Proceedings of the International Congress of Mathematicians (ICM 2018)*, 2019.
- Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, 2019.
- Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error, 2020.
- Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802, 2021.

## Appendix A. Proofs for Section 2

**Proof** [Proof of Theorem 10] To prove our first claim, we start by showing that there exists a policy  $\pi$  that achieves this value for each state. Let  $\pi$  be the policy which for any state  $s$  with assignment  $w \neq w^*$  chooses the action which decreases the hamming distance  $\text{dist}(w, w^*)$  by 1. Note that one such action always exists in our construction, since a satisfying assignment satisfies all clauses. Therefore, from a state  $s$  at level  $l$  with assignment  $w$ , we can reach a state with assignment  $w_1$  such that either (i)  $w_1$  is a satisfying assignment or (ii)  $w_1$  is at the last level. In both cases,

$$V^\pi(s) = \left(1 - \frac{l + \text{dist}(w, w_1) + \text{dist}(w_1, w^*)}{H + v}\right)^r = g(l, w)$$

Next, for any other policy  $\pi'$  that ends on state  $s'$  at level  $l'$  with assignment  $w'$  (i.e. either  $l' = H$  or  $w' = w^*$ ), we have

$$V^{\pi'}(s) = \left(1 - \frac{l' + \text{dist}(w', w^*)}{H + v}\right)^r \leq \left(1 - \frac{l + \text{dist}(w, w') + \text{dist}(w', w^*)}{H + v}\right)^r \leq g(l, w)$$

where the first inequality follows from  $l' - l \geq \text{dist}(w, w')$ . This proves our first claim about  $V^*$  i.e.  $V^*(s) = g(l, w)$ .

To prove our second claim, that  $V^*$  and  $Q^*$  can be written as a linear function of features  $\psi$ , we will show that  $V^*(s)$  can be written as a polynomial of degree at most  $r$  in  $w$  and  $w^*$ . To see why this is enough, we set  $\theta$  to be all monomials in  $w^*$  of degree at most  $r$ . That is, each coordinate of  $\theta$  corresponds to a multiset  $S \subset [v]$  of size  $|S| \leq r$ , and its value is  $\theta_S = \prod_{i \in S} w_i^*$ . We set  $\psi(s)$  to be the corresponding coefficients in the polynomial  $V^*$ . Then, we can write  $V^*(s) = \langle \theta, \psi(s) \rangle$ . Since, there are at most  $\sum_{i=0}^r v^i \leq 2v^r$  many coefficients we can set the feature dimension as  $d = 2v^r$ .

Finally, we prove that  $V^*(s)$  can be written as a polynomial of degree at most  $r$  in  $w$  and  $w^*$ . Firstly hamming distance  $\text{dist}(w, w^*)$  is linear in both  $w$  and  $w^*$  i.e.

$$\text{dist}(w, w^*) = \frac{v - \langle w, w^* \rangle}{2}$$

Our claim follows from noting that  $g(l, w)$  is a polynomial of degree  $r$  in  $\text{dist}(w, w^*)$ . Note that linear  $V^*$  implies linear  $Q^*$  in deterministic MDPs for  $\psi(s, a) = \psi(P(s, a))$ , since by definition, in MDPs with deterministic transition,  $Q^*(s, a) = V^*(P(s, a))$ .  $\blacksquare$

**Proof** [Proof of Theorem 8] For any  $q \geq 1$ , we set

$$r = 8q. \tag{2}$$

Note that  $q \geq 1$  implies  $r \geq 2$ . Therefore, to prove our proposition, we just need to show

$$d^q \leq v^{r^2/4} \tag{3}$$

$$d \cdot v^{r^2/4} \leq v^{8q+16q^2+1} \tag{4}$$

under this setting of  $d$  and  $r$ . Here the first equation bounds the time complexity of LINEAR-3-RL in terms of feature dimension  $d$  and the second equation bounds the time complexity of UNIQUE-3-SAT in terms of the number of variables  $v$ . Equation (3) is true as

$$v^{\frac{r^2}{4}} = (v^r)^{\frac{r}{4}} \geq d^{\frac{r}{8}} = d^q$$

where the first inequality follows from  $d \leq v^{2r}$  for large enough  $v$  and the last equality follows from Equation (2) above. Equation (4) holds since

$$d \cdot v^{r^2/4} = 2v^{r+r^2/4} = O(v^{8q+16q^2}),$$

where the first equality follows from  $d = 2v^r$  and the last equality follows from Equation (2) for large enough  $v$ .  $\blacksquare$

**Proof** [Proof of Theorem 9] This follows exactly as proof of Theorem 8. We set

$$r = \left\lceil \frac{\sqrt{v}}{\log v} \right\rceil.$$

We proceed to show that (i) the time complexity of LINEAR-3-RL can be bounded by  $v^{r^2/4}$  (ii)  $d \cdot v^{r^2/4}$ , which is the time complexity of UNIQUE-3-SAT if (i) holds, can be bounded by  $2^{O(v/\log v)}$ .

With this setting, the time complexity of LINEAR-3-RL simplifies to

$$d^{\frac{\log d}{32 \log \log d}} \leq v^{\frac{2r \log d}{32 \log \log d}} \leq v^{\frac{4r^2 \log v}{32 \log(r \log v)}} \leq v^{\frac{8r^2 \log v}{32 \log v}} = v^{\frac{r^2}{4}}$$

where the first and second inequality follows from  $v^r \leq d \leq v^{2r}$  and third inequality follows from our setting of  $r$ .

Similarly, the time complexity of UNIQUE-3-SAT simplifies to

$$d \cdot v^{r^2/4} = 2v^{r+\frac{r^2}{4}} \leq v^{r^2} \leq v^{\frac{4v}{\log^2 v}} = 2^{\frac{4v}{\log v}}$$

where the first equality follows from  $d = 2v^r$ , first inequality follows for  $r \geq 2$  and large enough  $v$  and the second inequality follows from our setting of  $r$  above.  $\blacksquare$

The following is generalization of reduction from UNIQUE-3-SAT to LINEAR-3-RL in Theorem 8.

**Proposition 14** *Suppose  $m \geq 0$  and  $q \geq 1$ . If LINEAR-3-RL with feature dimension  $d$  can be solved in time  $d^{q \cdot (\log d)^{m/(m+2)}}$  with error probability  $1/10$ , then UNIQUE-3-SAT with  $v$  variables can be solved in time  $v^{O(16^{m+2}q^{m+2}) \cdot (\log v)^m}$  with error probability  $1/8$ .*

**Proof** For any  $m \geq 0$  and  $q \geq 1$ , we set

$$r = \left\lceil \sqrt{(16q)^{m+2} \log^m v} \right\rceil. \quad (5)$$

Note that  $m \geq 0$  and  $q \geq 1$  implies  $r \geq 2$ . Therefore, to prove our claim, we just need to show the following equations hold for our setting of  $d$  and  $r$ :

$$v^{r^2/4} \geq d^{q \cdot \log^{\frac{m}{m+2}} d} \quad (6)$$

$$d \cdot v^{r^2/4} = v^{O((16q)^{m+2}) \log^m v} \quad (7)$$

Here the first equation bounds the time complexity of LINEAR-3-RL in terms of feature dimension  $d$  and the second equation bounds the time complexity of UNIQUE-3-SAT in terms of the number of variables  $v$ .

**Proof of Equation (6):** To prove the first inequality, we lower bound  $r$  in terms of feature dimension  $d$  as

$$r \geq 8q (\log d)^{m/(m+2)}. \quad (8)$$

which can be proved by lower bounding  $r^{m+2}$  as follows

$$\begin{aligned} r^{m+2} &= r^2 \cdot r^m \geq (16q)^{m+2} \log^m v \cdot r^m \\ &= (16q)^{m+2} \log^m(v^r) \geq (16q)^{m+2} \log^m(\sqrt{d}) \geq (8q)^{m+2} \log^m d \end{aligned}$$

where the first inequality follows from our setting of  $r$  and the second inequality follows from  $d \leq 2v^r \leq v^{2r}$  for  $r > 0$  and large enough  $v$ . Substituting the lower bound in  $v^{r^2/4}$ , we can write the time complexity of LINEAR-3-RL in terms of feature dimension  $d$  as

$$v^{\frac{r^2}{4}} = (v^r)^{\frac{r}{4}} \geq d^{\frac{r}{8}} \geq d^{q \cdot \log \frac{m}{m+2}} d$$

where the first inequality follows again from  $d \leq v^{2r}$  and the second inequality follows from Equation (8) above.

**Proof of Equation (7):** The second equation follows by substituting our setting of  $r$  (Equation (2)) in  $d \cdot v^{r^2/4}$ ,

$$d \cdot v^{r^2/4} = 2v^{r+r^2/4} = v^{O((16q)^{m+2}) \log^m v}$$

where the first equality follows from  $d = 2v^r$ . ■

## Appendix B. Proofs for Section 3

**Proof [Proof of Theorem 13]** The proof for the value function of non-intermediate state is identical to that in the 3-action MDP. We proceed to argue the second claim. For an intermediate state, the value function will be identical to its parent if the two actions available include a wrong bit that ought to be flipped in the optimal assignment  $w^*$ . Otherwise, no matter what action the agent takes, it will reach a non-intermediate state whose depth is  $l + 1$  and hamming distance is  $\text{dist}(w, w^*) + 1$ . Compared to the value function of its parent, such intermediate state will have an extra 2 in the numerator. We encode the situation with the indicator term  $\mathbb{1}\{w_{i_1} = w_{i_1}^*\} \cdot \mathbb{1}\{w_{i_2} = w_{i_2}^*\}$ . This then gives the value function for these intermediate states.

Lastly, like in the proof of Theorem 10, it suffices to argue the value function is a degree  $2r$  polynomial in  $w$  and  $w^*$ . This is by noticing that (i)  $\text{dist}(w, w^*)$  is linear in  $w$  and  $w^*$ ; and (ii)  $\mathbb{1}\{w_{i_1} = w_{i_1}^*\} \cdot \mathbb{1}\{w_{i_2} = w_{i_2}^*\}$  is quadratic in  $w$  and  $w^*$  i.e.

$$\mathbb{1}\{w_{i_1} = w_{i_1}^*\} \cdot \mathbb{1}\{w_{i_2} = w_{i_2}^*\} = \frac{1}{4} \cdot (1 - w_{i_1} \cdot w_{i_1}^*) (1 - w_{i_2} \cdot w_{i_2}^*).$$

Thus, the value function is overall a polynomial of degree  $2r$  in both  $w$  and  $w^*$ . This proves the last claim. ■

**Proof [Proof of Theorem 2]** Fix  $q \geq 1$ , we set  $r = 12q$ . Under this setting, we have

$$d^q \leq d^{r/12} \leq (v^{3r})^{r/12} = v^{r^2/4},$$

where the first inequality follows from the setting of  $r$  and the second inequality follows from  $d \leq v^{3r}$ . The reduction then allows us to upper bound the complexity of UNIQUE-3-SAT by

$$d \cdot v^{r^2/4} \leq v^{r^2/4+3r} = v^{O(q^2)},$$

where the first inequality again follows from  $d \leq v^{3r}$ . ■

**Proof** [Proof of Theorem 6] The proof follows similarly as proof of Theorem 9. The only difference is that since  $d$  is now bounded by  $v^{3r}$  instead of  $v^{2r}$ , we need the runtime of LINEAR-2-RL in the assumption to also have a different constant in the exponent i.e.  $d^{\log d / (72 \log \log d)}$ . ■