
Going Beyond Linear RL: Sample Efficient Neural Function Approximation

Baihe Huang^{1*}, Kaixuan Huang^{2*}, Sham M. Kakade^{3,4*}, Jason D. Lee^{2*}
Qi Lei^{2*}, Runzhe Wang^{2*}, Jiaqi Yang^{5*}

¹Peking University ²Princeton University ³Harvard University
⁴Microsoft Research ⁵Tsinghua University

Abstract

Deep Reinforcement Learning (RL) powered by neural net approximation of the Q function has had enormous empirical success. While the theory of RL has traditionally focused on linear function approximation (or eluder dimension) approaches, little is known about nonlinear RL with neural net approximations of the Q functions. This is the focus of this work, where we study function approximation with two-layer neural networks (considering both ReLU and polynomial activation functions). Our first result is a computationally and statistically efficient algorithm in the generative model setting under completeness for two-layer neural networks. Our second result considers this setting but under only realizability of the neural net function class. Here, assuming deterministic dynamics, the sample complexity scales linearly in the algebraic dimension. In all cases, our results significantly improve upon what can be attained with linear (or eluder dimension) methods.

1 Introduction

In reinforcement learning (RL), an agent aims to learn the optimal decision-making rule by interacting with an unknown environment [71]. Deep Reinforcement Learning, empowered by deep neural networks [48, 32], has achieved tremendous success in various real-world applications, such as Go [68], Atari [54], Dota2 [7], Texas Holdém poker [56], and autonomous driving [66]. Those modern RL applications are characterized by large state-action spaces, and the empirical success of deep RL corroborates the observation that deep neural networks can extrapolate well across state-action spaces [35, 55, 50].

Although in practice non-linear function approximation scheme is prevalent, theoretical understandings of the sample complexity of RL mainly focus on tabular or linear function approximation settings [69, 38, 5, 41, 63, 88, 1, 43, 44, 77]. These results rely on finite state space or exact linear approximations. Recently, sample efficient algorithms under non-linear function approximation settings are proposed [82, 11, 20, 12, 51, 72, 13, 91, 86]. Those algorithms are based on Bellman rank [40], eluder dimension [64], neural tangent kernel [37, 4, 14, 92], or sequential Rademacher complexity [60, 61]. Yet, the understanding of how deep RL learns and generalizes in large state spaces is far from complete. While the aforementioned works study function approximation structures that possess the nice properties of linear models, such as low information gain and low eluder dimensions, the highly non-linear nature of neural networks renders challenges on their applicability to deep RL. For one thing, recent wisdoms in deep learning theory cast doubt on the ability of neural tangent

*Alphabetical order. Correspondence to: Baihe Huang, baihehuang@pku.edu.cn, Kaixuan Huang, kaixuanh@princeton.edu, Sham M. Kakade, sham@cs.washington.edu, Jason D. Lee, JasonDL@princeton.edu, Qi Lei, qilei@princeton.edu, Runzhe Wang, runzhe@princeton.edu, and Jiaqi Yang, yangjq17@gmail.com.

kernel and random features to model the actual neural networks. Indeed, the neural tangent kernel approximately reduces neural networks to linear models, but the RKHS norm of neural networks is exponential [87]. Moreover, it remains unclear what neural networks models have low eluder dimensions. For example, recent work [13] shows that two layer neural networks have exponential eluder dimension in the dimension of features. Thus, the mismatch between the empirical success of deep RL and its theoretical understanding remains significant, which yields the following important question:

What are the structural properties that allow sample-efficient algorithms for RL with neural network function approximation?

Recent work in RL suggests that learning RL with neural networks function approximation is exponentially hard [13, 49, 52]. In this paper, however, we advance the understanding of the above question by displaying several structural properties that allow efficient RL algorithms with neural function approximations. We consider several value function approximation models that possess high information gain and high eluder dimension. Specifically, we study two structures, namely two-layer neural networks and structured polynomials (i.e. two-layer neural networks with polynomial activation functions), under two RL settings, namely RL with simulator model and online RL. In the simulator (generative model) setting [45, 67], the agent can simulate the MDP at any state-action pair. In online RL, the agent can only start at an initial state and interact with the MDP step by step. The goal in both settings is to find a near-optimal policy while minimizing the number of samples used.

We obtain the following results. For the simulator setting, we propose sample-efficient algorithms for RL with two-layer neural network function approximation. Under either policy completeness, Bellman completeness, or gap conditions, our method provably learns near-optimal policy with polynomial sample complexities. For online RL, we provide sample-efficient algorithms for RL with structured polynomial function approximation. When the transition is deterministic, we also present sample-efficient algorithms under only the realizability assumption [18, 78]. Our main techniques are based on neural network recovery [90, 39, 28], and algebraic geometry [53, 65, 8, 76].

1.1 Summary of our results

Our main results in different settings are summarized in Table 1. We consider two-layer neural networks $f(x) = \langle v, \sigma(Wx) \rangle$ (where σ is ReLU activation) and rank k polynomials (see Example 4.3).

Table 1: Baselines and our main results for the sample complexity to find an ϵ -optimal policy.

	rank k polynomial		Neural Net of Width k		
	Sim. + Det. (R)	Onl. + Det. (R)	Sim. + Det. (R)	Sim. + Gap. (R)	Sim. + Stoch. (C)
Baseline	$O(d^p)$	$O(d^p)$	$O(d^{\text{poly}(1/\epsilon)})$ (*)	$O(d^{\text{poly}(1/\epsilon)})$	$O(d^{\text{poly}(1/\epsilon)})$
Our results	$O(dk)$	$O(dk)$	$\tilde{O}(\text{poly}(d) \cdot \exp(k))$	$\tilde{O}(\text{poly}(d, k))$	$\tilde{O}(\text{poly}(d, k)/\epsilon^2)$

We make the following elaborations on Table 1.

- For simplicity, we display only the dependence on the feature dimension d , network width or polynomial rank k , precision ϵ , and degree p (of polynomials).
- In the table *Sim.* denotes simulator model, *Onl.* denotes online RL, *Det.* denotes deterministic transitions, *Stoch.* denotes stochastic transitions, *Gap.* denotes gap condition, *(R)* denotes realizability assumption only, and *(C)* denotes completeness assumption (either policy complete or Bellman complete) together with realizability assumption.
- We apply [21] for the deterministic transition baseline, and apply [19] for the stochastic transition baseline. We are unaware of any methods that can directly learn MDP with neural network value function approximation².
- In polynomial case, the baseline first vectorizes the tensor $\begin{pmatrix} 1 \\ x \end{pmatrix}^{\otimes p}$ into a vector in $\mathbb{R}^{(d+1)^p}$ and then performs on this vector. In the neural network case, the baseline uses a polynomial of degree $1/\epsilon$ to approximate the neural network with precision ϵ and then vectorizes the

polynomial into a vector in $\mathbb{R}^{d^{\text{poly}(1/\epsilon)}}$. The baseline method for realizable model (denoted by $(*)$) needs a further gap assumption of $\text{gap} \geq d^{\text{poly}(1/\epsilon)}\epsilon$ to avoid the approximation error from escalating [21]; note for small ϵ this condition never holds but we include it in the table for the sake of comparison.

- In rank k polynomial case, our result $O(dk)$ in simulator model can be found in Theorem 4.7 and our result $O(dk)$ in online RL model can be found in Theorem 4.8. These results only require a realizability assumption. Efficient explorations are guaranteed by algebraic-geometric arguments. In neural network model, our result $\tilde{O}(\text{poly}(d) \cdot \exp(k))$ in simulator model can be found in Theorem 3.4. This result also only relies on the realizability assumption. For stochastic transitions, our result $\tilde{O}(\text{poly}(d, k)/\epsilon^2)$ works for either policy complete or Bellman complete settings, as in Theorem 3.5 and Theorem 3.6 respectively. The $\tilde{O}(\text{poly}(d, k))$ result for gap condition can be found in Theorem 3.8.

1.2 Related Work

Linear Function Approximation. RL with linear function approximation has been widely studied under various settings, including linear MDP and linear mixture MDP [44, 89, 85]. While these papers have proved efficient regret and sample complexity bounds, their analyses relied heavily on two techniques: they used the confidence ellipsoid to quantify the uncertainty, and they used the elliptical potential lemma to bound the total uncertainty [2]. These two techniques were integral to their analyses but are so restrictive that they generally do not extend to nonlinear cases.

Eluder Dimension. [62, 59] proposed eluder dimension, a complexity measure of the function space, and proved regret and sample complexity bounds that scaled with the eluder dimension, for bandits and reinforcement learning [73, 42]. They also showed that the eluder dimension is small in several settings, including generalized linear models and LQR. However, as shown in [13], the eluder dimension could be exponentially large even with a single ReLU neuron, which suggested the eluder dimension would face difficulty in dealing with neural network cases. The eluder dimension is only known to give non-trivial bounds for linear function classes and monotone functions of linear function classes. For structured polynomial classes, the eluder dimension simply embeds into an ambient linear space of dimension d^p , where d is the dimension, and p is the degree. This parallels the lower bounds in linearization / neural tangent kernel (NTK) works [79, 30, 3], which show that linearization also incurs a similarly large penalty of d^p sample complexity, and more advanced algorithm design is need to circumvent linearization[6, 9, 22, 83, 26, 58, 29, 57, 34, 75, 10].

Bellman Rank and Completeness. [40, 70] studied RL with general function approximation. They used Bellman rank to measure the error of the function class under the Bellman operator and gave proved bounds in the term of it. Recently, [16] propose bilinear rank and encompass more function approximation models. However, it is hard to bound either the Bellman rank or the bilinear rank for neural nets. Therefore, their results are not known to include the neural network approximation setting. Another line of work shows that exponential sample complexity is unavoidable even with good representations [19, 80, 78], which implies the realizability assumption alone might be insufficient for function approximations.

2 Preliminaries

An episodic Markov Decision Process (MDP) is defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$ where \mathcal{S} is the state space, \mathcal{A} is the action set, H is the number of time steps in each episode, \mathbb{P} is the transition kernel and r is the reward function. In each episode the agent starts at a fixed initial state s_1 and at each time step $h \in [H]$ it takes action a_h , receives reward $r_h(s_h, a_h)$ and transits to $s_{h+1} \sim \mathbb{P}(\cdot | s_h, a_h)$.

A deterministic policy π is a length- H sequence of functions $\pi = \{\pi_h : \mathcal{S} \mapsto \mathcal{A}\}_{h=1}^H$. Given a policy π , we define the value function $V_h^\pi(s)$ as the expected sum of reward under policy π starting from

²Prior work on neural function approximation has focused on neural tangent kernels, which would require $d^{\text{poly}(1/\epsilon)}$ to approximate a two-layer network [31].

$s_h = s$:

$$V_h^\pi(s) := \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) | s_h = s \right]$$

and we define the Q function $Q_h^\pi(s, a)$ as the the expected sum of reward taking action a in state $s_h = s$ and then following π :

$$Q_h^\pi(s, a) := \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) | s_h = s, a_h = a \right].$$

The Bellman operator \mathcal{T}_h applied to Q-function Q_{h+1} is defined as follow

$$\mathcal{T}_h(Q_{h+1})(s, a) := r_h(s, a) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} [\max_{a'} Q_{h+1}(s', a')].$$

There exists an optimal policy π^* that gives the optimal value function for all states, i.e. $V_h^{\pi^*}(s) = \sup_{\pi} V_h^\pi(s)$ for all $h \in [H]$ and $s \in \mathcal{S}$. For notational simplicity we abbreviate V^{π^*} as V^* and correspondingly Q^{π^*} as Q^* . Therefore Q^* satisfies the following Bellman optimality equations for all $s \in \mathcal{S}$, $a \in \mathcal{A}$ and $h \in [H]$:

$$Q_h^*(s, a) = \mathcal{T}_h(Q_{h+1}^*)(s, a).$$

The goal is to find a policy π that is ϵ -optimal in the sense that $V_1^*(s_1) - V_1^\pi(s_1) \leq \epsilon$, within a small number of samples. We consider two query models of interacting with MDP:

- In the simulator model ([45], [67]), the agent interacts with a black-box that simulates the MDP. At each time step $h \in [H]$, the agent can start at a state-action pair (s, a) and interact with the black box by executing some policy π chosen by the agent.
- In online RL, the agent can only start at the initial state and interact with the MDP by using a policy and observing the rewards and the next states. In each episode k , the agent proposes a policy π^k based on all history up to episode $k - 1$ and executes π^k to generate a single trajectory $\{s_h^k, a_h^k\}_{h=1}^H$ with $a_h^k = \pi_h^k(s_h^k)$ and $s_{h+1}^k \sim \mathbb{P}(\cdot | s_h^k, a_h^k)$.

2.1 Function approximation

In reinforcement learning with value function approximation, the learner is given a function class $\mathcal{F} = \mathcal{F}_1 \times \dots \times \mathcal{F}_H$, where $\mathcal{F}_h \subset \{f : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]\}$ is a set of candidate functions to approximate Q^* . The following assumption is commonly adopted in the literature [43, 74, 42, 17].

Assumption 2.1 (Realizability). $Q_h^* \in \mathcal{F}_h$ for all $h \in [H]$.

The function approximation is equipped with feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \{u \in \mathbb{R}^d : \|u\|_2 \leq B_\phi\}$ that is known to the agent. We focus the continuous action setting (e.g. in control and robotics problems) and make the following regularity assumption on the feature function ϕ .

Assumption 2.2 (Bounded features). Assume $\phi(s, a) \leq B_\phi, \forall (s, a) \in \mathcal{S} \times \mathcal{A}$.

Notation For any vector $x \in \mathbb{R}^d$, let $x_{\max} := \max_{i \in [d]} x_i$ and $x_{\min} := \min_{i \in [d]} x_i$. Let $s_i(\cdot)$ denote the i -th singular value, $s_{\min}(\cdot)$ denotes the minimum eigenvalue and $s_{\max}(\cdot)$ denotes the maximum eigenvalue. The conditional number is defined by $\kappa(\cdot) := s_{\max}(\cdot)/s_{\min}(\cdot)$. We use \otimes to denote Kronecker product and \circ to denote Hadamard product. For a given integer H , we use $[H]$ to denote the set $\{1, 2, \dots, H\}$. For a function $f : \mathfrak{X} \mapsto \mathfrak{Y}$, we use $f^{-1}(y) := \{x \in \mathfrak{X} : f(x) = y\}$ to denote the preimage of $y \in \mathfrak{Y}$. We use the shorthand $x \lesssim y$ ($x \gtrsim y$) to indicate $x \leq O(y)$ ($x \geq \Omega(y)$).

3 Neural Network Function Approximation

In this section we show sample-efficient algorithms with neural network function approximations. The function class of interest is given in the following definition. More general neural network class is discussed in Appendix B.5.

Definition 3.1 (Neural Network Function Class). We use \mathcal{F}_{NN} to denote the function class of $f(\phi(s, a)) : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ where $f(x) = \langle v, \sigma(Wx) \rangle : \mathbb{R}^d \mapsto \mathbb{R}$ is a two-layer neural network where σ is ReLU, $\|W\|_F \leq B_W$, $v \in \{\pm 1\}^k$, $\prod_{i=1}^k s_i(W)/s_{\min}(W) \leq \lambda$, $s_{\max}(W)/s_{\min}(W) \leq \kappa$ and $k \leq d$. Here $\phi : \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}^d$ is a known feature map whose image contains a ball $\{u \in \mathbb{R}^d : \|u\|_2 \leq \delta_\phi\}$ with $\delta_\phi \geq d \cdot \text{polylog}(d)$.³

We introduce the following completeness properties in the setting of value function approximations. Along with Assumption 2.1, they are commonly adopted in the literature .

Definition 3.2 (Policy complete). Given MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, H)$, function class $\mathcal{F}_h : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, $h \in [H]$ is called policy complete if for all π and $h \in [H]$, $Q_h^\pi \in \mathcal{F}_h$.

Definition 3.3 (Bellman complete). Given MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, H)$, function class $\mathcal{F}_h : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, $h \in [H]$ is called Bellman complete if for all $h \in [H]$ and $Q_{h+1} \in \mathcal{F}_{h+1}$, $\mathcal{T}_h(Q_{h+1}) \in \mathcal{F}_h$.

3.1 Warmup: Realizable Q^* with deterministic transition

We start by considering the case when the transition kernel is deterministic. In this case only Assumption 2.1 is required for the expressiveness of neural network function approximations. Algorithm 1 learns optimal policy from time step H to 1. Suppose we have learned policies π_{h+1}, \dots, π_H at level h and they are exactly the optimal policies. We first explore features $\phi(s_h^i, a_h^i)$ over a standard Gaussian distribution, and if $\|\phi(s_h^i, a_h^i)\|_2 \geq \delta_\phi$ then we simply skip this trial. Recall that $\delta_\phi \geq d \cdot \text{poly log}(d)$, so with high probability (w.r.t d) almost all feature samples will be explored. We next construct an estimate \widehat{Q}_h^i of $Q^*(s_h^i, a_h^i)$ by collecting cumulative rewards using π_{h+1}, \dots, π_H as the roll-out. Since the transition is deterministic, $\widehat{Q}_h^i = Q^*(s_h^i, a_h^i)$ for all explored samples (s_h^i, a_h^i) . Recall that Q_h^* is a two-layer neural network, we can now recover its parameters in Line 12 exactly by invoking techniques in neural network optimization (see, e.g. [39, 27, 90]). Details of this step can be found in Appendix B.5, where the method is mainly based on [90]. This means the reconstructed \widehat{Q}_h in Line 13 is precisely Q^* , and the algorithm can thus find optimal policy π_h^* in the h -th level.

Algorithm 1 Learning realizable Q^* with deterministic transition

```

1: for  $h = H, \dots, 1$  do
2:   Sample  $x_h^i, i \in [n]$  from standard Gaussian  $\mathcal{N}(0, I_d)$ 
3:   for  $i \in [n]$  do
4:     if  $\|x_h^i\| \leq \delta_\phi$  then
5:       Find  $(s_h^i, a_h^i) \in \phi^{-1}(x_h^i)$  and locate the state  $s_h^i$  in the generative model
6:       Pull action  $a_h^i$  and use  $\pi_{h+1}, \dots, \pi_H$  as the roll-out to collect rewards  $r_h^{(i)}, \dots, r_H^{(i)}$ 
7:       Construct estimation
           
$$\widehat{Q}_h^i \leftarrow r_h^{(i)} + \dots + r_H^{(i)}$$

8:     else
9:       Let  $\widehat{Q}_h^i \leftarrow 0$ .
10:    end if
11:  end for
12:  Compute  $(v_h, W_h) \leftarrow \text{NEURALNETRECOVERY}(\{(x_h^i, \widehat{Q}_h^i) : i \in [n]\})$ 
13:  Set  $\widehat{Q}_h(s, a) \leftarrow v_h^\top \sigma(W_h \phi(s, a))$ 
14:  Let  $\pi_h(s) \leftarrow \arg \max_{a \in \mathcal{S}} \widehat{Q}_h(s, a)$ 
15: end for
16: Return  $\pi_1, \dots, \pi_H$ 

```

Theorem 3.4. (Informal) If $n \geq d \cdot \text{poly}(\kappa, k, \lambda, \log d, B_W, B_\phi, H)$, then with high probability Algorithm 1 learns the optimal policy.

The formal statement and complete proof are deferred to the Appendix B.1. The main idea of exact neural network recovery can be summarized in the following. We first use method of moments to

³Here the δ_ϕ is chosen only for simplicity. In general this assumption can be relaxed to that the image of ϕ contains an arbitrary dense ball near the origin, since one can always rescale the feature mapping in the neural function approximation.

find a ‘rough’ parameter recovery. If this ‘rough’ recovery is sufficiently close to the true parameter, the empirical ℓ_2 loss function is locally strongly convex and there is unique global minimum. Then we can apply gradient descent to find this global minimum which is exactly the true parameter.

3.2 Policy complete neural function approximation

Now we consider general stochastic transitions. Difficulties arise in this scenario due to noises in the estimation of Q-functions. In the presence of model misspecification, these noises cause estimation errors to amplify through levels and require samples to be exponential in H . In this section, we show that neural network function approximation is still learnable, assuming the function class \mathcal{F}_{NN} is policy complete with regard to MDP \mathcal{M} . Thus for all $\pi \in \Pi$, we can denote $Q_h^\pi(s, a) = \langle v^\pi, \sigma(W^\pi \phi(s, a)) \rangle$.

Algorithm 2 Learn policy complete NN with simulator.

```

1: for  $h = H, \dots, 1$  do
2:   Sample  $x_h^i, i \in [n]$  from standard Gaussian  $\mathcal{N}(0, I_d)$ 
3:   for  $i \in [n]$  do
4:     if  $\|x_h^i\| \leq \delta_\phi$  then
5:       Find  $(s_h^i, a_h^i) \in \phi^{-1}(x_h^i)$  and locate the state  $s_h^i$  in the generative model
6:       Pull action  $a_h^i$  and use  $\pi_{h+1}, \dots, \pi_H$  as the roll-out to collect rewards  $r_h^{(i)}, \dots, r_H^{(i)}$ 
7:       Construct unbiased estimation of  $Q_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i)$ 
           
$$\widehat{Q}_h^i \leftarrow r_h^{(i)} + \dots + r_H^{(i)}$$

8:     else
9:       Let  $\widehat{Q}_h^i \leftarrow 0$ .
10:    end if
11:  end for
12:  Compute  $(v_h, W_h) \leftarrow \text{NEURALNETNOISYRECOVERY}(\{(x_h^i, \widehat{Q}_h^i) : i \in [n]\})$ 
13:  Set  $\widehat{Q}_h(s, a) \leftarrow v_h^\top \sigma(W_h \phi(s, a))$ 
14:  Let  $\pi_h(s) \leftarrow \arg \max_{a \in \mathcal{S}} \widehat{Q}_h(s, a)$ 
15: end for
16: Return  $\pi_1, \dots, \pi_H$ 

```

Algorithm 2 learns policy from level $H, H-1, \dots, 1$. In level h , the algorithm has learned policy π_{h+1}, \dots, π_H that is only sub-optimal by $(H-h)\epsilon/H$. Then it explores features $\phi(s, a)$ from $\mathcal{N}(0, I_d)$. The algorithm then queries (s, a) and uses learned policy π_{h+1}, \dots, π_H as roll out to collect an unbiased estimate of the Q-function $Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a)$. Since $Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a) \in \mathcal{F}_{NN}$ is a two-layer neural network, it can then be recovered from samples. Details of this step can be found in Appendix B.5, where the methods are mainly based on [90]. The algorithm then reconstructs this Q-function and finds its optimal policy π_h .

Theorem 3.5. (Informal) Fix ϵ, t , if $n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ Algorithm 2 returns an ϵ -optimal policy π .

The formal statement and complete proof are deferred to Appendix B.2. Notice that unlike the case of Theorem 3.4, the sample complexity does not depend on λ , thus avoiding the potential exponential dependence in k .

The main idea of the proof is that at each time step a neural network surrogate of Q^* can be constructed by the policy already learned. Suppose we have learned π_{h+1}, \dots, π_H in level h , then from policy completeness $Q_h^{\pi_{h+1}, \dots, \pi_H}$ belongs to \mathcal{F}_{NN} and we can interact with the simulator to obtain its estimate \widehat{Q}_h . If $\|\widehat{Q}_h - Q_h^{\pi_{h+1}, \dots, \pi_H}\|_\infty$ is small, the optimistic planning based on \widehat{Q}_h is not far from the optimal policy of $Q_h^{\pi_{h+1}, \dots, \pi_H}$. Therefore the errors can be decoupled into the errors in recovering $Q_h^{\pi_{h+1}, \dots, \pi_H}$ and the suboptimality of $Q_h^{\pi_{h+1}, \dots, \pi_H}$, which depends on level $h+1$. This reasoning can then be recursively performed to level H , and thus we can bound the suboptimality of π_h .

3.3 Bellman complete neural function approximation

In addition to policy completeness, we show that neural network function approximation can also learn efficiently under the setting where the function class \mathcal{F}_{NN} is Bellman complete with regard to MDP \mathcal{M} . Specifically, for $Q_{h+1} \in \mathcal{F}_{h+1}$, there are $v^{Q_{h+1}}$ and $W^{Q_{h+1}}$ such that $\mathcal{T}_h(Q_{h+1})(s, a) = \langle v^{Q_{h+1}}, \sigma(W^{Q_{h+1}}\phi(s, a)) \rangle$.

Algorithm 3 is similar to the algorithm in previous section. Suppose in level h , the algorithm has constructed the Q-function $\widehat{Q}_{h+1}(s, a) = v_{h+1}^\top \sigma(W_{h+1}\phi(s, a))$ that is $(H - h)\epsilon/H$ -close to the optimal Q_{h+1}^* . It then recovers weights v_h, W_h from $\mathcal{T}_h(\widehat{Q}_{h+1})(s, a) = \langle v^{\widehat{Q}_{h+1}}, \sigma(W^{\widehat{Q}_{h+1}}\phi(s, a)) \rangle$, using unbiased estimates $r_h(s_h^i, a_h^i) + \widehat{V}_{h+1}(s_{h+1}^i)$. The Q-function $\widehat{Q}_h(s, a) = v_h^\top \sigma(W_h\phi(s, a))$ reconstructed from weights v_h, W_h is thus $(H - h + 1)\epsilon/H$ -close to the Q_h^* .

Algorithm 3 Learn Bellman complete NN with simulator.

```

1: for  $h = H, \dots, 1$  do
2:   Sample  $x_h^i, i \in [n]$  from standard Gaussian  $\mathcal{N}(0, I_d)$ 
3:   for  $i \in [n]$  do
4:     if  $\|x_h^i\| \leq \delta_\phi$  then
5:       Find  $(s_h^i, a_h^i) \in \phi^{-1}(x_h^i)$  and locate the state  $s_h^i$  in the generative model
6:       Pull action  $a_h^i$  and observe  $r_h(s_h^i, a_h^i), s_{h+1}^i$ 
7:       Construct unbiased estimation of  $\mathcal{T}_h(\widehat{Q}_{h+1})(s_h^i, a_h^i)$ 
           
$$\widehat{Q}_h^i \leftarrow r_h(s_h^i, a_h^i) + \widehat{V}_{h+1}(s_{h+1}^i)$$

8:     else
9:       Let  $\widehat{Q}_h^i \leftarrow 0$ .
10:    end if
11:  end for
12:  Compute  $(v_h, W_h) \leftarrow \text{NEURALNETNOISYRECOVERY}(\{(x_h^i, \widehat{Q}_h^i) : i \in [n]\})$ 
13:  Set  $\widehat{Q}_h(s, a) \leftarrow v_h^\top \sigma(W_h\phi(s, a))$  and  $\widehat{V}_h \leftarrow \max_{a \in \mathcal{A}} \widehat{Q}_h(s, a)$ 
14:  Let  $\pi_h(s) \leftarrow \arg \max_{a \in \mathcal{A}} \widehat{Q}_h(s, a)$ 
15: end for
16: Return  $\pi_1, \dots, \pi_H$ 

```

Theorem 3.6. (Informal) Fix ϵ, t , if $n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ Algorithm 3 returns an ϵ -optimal policy π .

Due to Bellman completeness, the error of estimation \widehat{Q}_h can be controlled recursively. In fact, we can show $\|\widehat{Q}_h - Q^*(s, a)\|_\infty$ is small by induction. The formal statement and detailed proof are deferred to Appendix B.3. Similar to Theorem 3.5, the sample complexity does not explicitly depend on λ , thus avoiding potentially exponential dependence in k .

3.4 Realizable Q^* with optimality gap

In this section we consider MDPs where there is a non-zero gap between the optimal policy and any other ones. This concept, known as optimality gap, is widely used in reinforcement learning and bandit literature [20, 19, 21].

Definition 3.7. The optimality gap is defined as

$$\rho = \inf_{a: Q^*(s, a) \neq V^*(s)} V^*(s) - Q^*(s, a).$$

We show that in the presence positive optimality gap, there exists an algorithm that can learn the optimal policy with polynomial samples even without the completeness assumptions. Intuitively, this is because one only needs to recover the neural network up to precision $\rho/4$ in order to make sure the greedy policy is identical to the optimal one. The formal statement and proof are deferred to Appendix B.4.

Theorem 3.8. (Informal) Fix $t \in (0, 1)$, if $n = \frac{d}{\rho^2} \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ there exists an algorithm that returns the optimal policy π^* .

Remark 3.9. In all aforementioned methods, there are two key components that allow efficient learning. First, the exploration is conducted in a way that guarantees an ℓ_∞ recovery of candidate functions. By ℓ_∞ recovery we mean the algorithm recovers a candidate Q -function in this class deviating from the target function Q^* by at most ϵ uniformly for all state-action pairs in the domain of interest. This notion of learning guarantee has received study in active learning [33, 46] and recently gain interest in contextual bandits [25]. Second, the agent constructs unbiased estimators of certain approximations to Q^* that lie in the neural function approximation class. This allows the recovery error to decouple linearly across time steps, which is made possible in several well-posed MDP instances, such as deterministic MDPs, MDPs with completeness assumptions, and MDPs with gap conditions. In principle, we note that provably efficient RL algorithms with general function approximation is possible as long as the above two components are present. We will see in the next section another example of learning RL with highly non-convex function approximations, where the function class of interest, admissible polynomial families, also allows for exploration schemes to achieve ℓ_∞ recovery.

4 Polynomial Realizability

In this section, we study the sample complexity to learn deterministic MDPs under polynomial realizability. We identify sufficient and necessary conditions for efficiently learning the MDPs for two different settings — the generative model setting and the online RL setting. Specifically, we show that if the image of the feature map $\phi_h(s_h, a_h)$ satisfies some positive measure conditions, then by random exploring, we can identify the optimal policy with samples linear in the algebraic dimension of the underlying polynomial class. We also provide a lower bound example showing the separation between the two settings.

Next, we introduce the notion of **Admissible Polynomial Families**, which are the families of structured polynomials that enable efficient learning.

Definition 4.1 (Admissible Polynomial Families). For $x \in \mathbb{R}^d$, denote $\tilde{x} = [1, x^\top]^\top$. Let $\mathcal{X} := \{\tilde{x}^{\otimes p} : x \in \mathbb{R}^d\}$. For any algebraic variety \mathcal{V} , we define $\mathcal{F}_\mathcal{V} := \{f_\Theta(x) = \langle \Theta, \tilde{x}^{\otimes p} \rangle : \Theta \in \mathcal{V}\}$ as the polynomial family parameterized by $\Theta \in \mathcal{V}$. We say $\mathcal{F}_\mathcal{V}$ is admissible⁴ w.r.t. \mathcal{X} , if for any $\Theta \in \mathcal{V}$, $\dim(\mathcal{X} \cap \{X \in \mathcal{X} : \langle X, \Theta \rangle = 0\}) < \dim(\mathcal{X}) = d$. We define the dimension D of the family to be the dimension of \mathcal{V} .

The following theorem shows that to learn an admissible polynomial family, the sample complexity only scales with the algebraic dimension of the polynomial family.

Theorem 4.2 ([36]). Consider the polynomial family $\mathcal{F}_\mathcal{V}$ of dimension D . For $n \geq 2D$, there exists a Lebesgue-measure zero set $N \in \mathbb{R}^d \times \dots \times \mathbb{R}^d$, such that if $(x_1, \dots, x_n) \notin N$, for any y_i , there is a unique f (or no such f) to the system of equations $y_i = f(x_i)$ for $f \in \mathcal{F}_\mathcal{V}$.

We give two important examples of admissible polynomial families with low dimension.

Example 4.3. (Low-rank Polynomial of rank k) The function $f \in \mathcal{F}_\mathcal{V}$ is a polynomial with k terms, that is

$$F(x) = \sum_{i=1}^k \lambda_i \langle v_i, x \rangle^{p_i},$$

where $p = \max\{p_i\}$. The dimension of this family is upper bounded by $D \leq dk$. Neural network with monomial/polynomial activation functions are low-rank polynomials.

Example 4.4. The function $f \in \mathcal{F}_\mathcal{V}$ is of the form $f(x) = q(Ux)$, where $U \in \mathbb{R}^{k \times d}$ and q is a degree p polynomial. The polynomial q and matrix U are unknown. The dimension of this family is upper bounded by $D \leq d(k+1)^p$.

Next, we introduce the notion of positive measure.

⁴Admissible means the dimension of \mathcal{X} decreases by one when there is an additional linear constraint $\langle \Theta, X \rangle = 0$

Definition 4.5. We say a measurable set $E \in \mathbb{R}^d$ is of positive measure if $\mu(E) > 0$, where μ is the standard Lebesgue measure on \mathbb{R}^d .

If a measurable set E satisfies $\mu(E) > 0$, then there exists a procedure to draw samples from E , such that for any $N \subset \mathbb{R}^d$ of Lebesgue-measure zero, the probability that the sample falls in N is zero. In fact, the sampling probability can be given by $\mathbb{P}_{x \in \mathcal{N}(0, I_d)}(\cdot | x \in E)$. The intuition behind its definition is that for all admissible polynomial families, the set of (x_1, \dots, x_n) with "redundant information" about learning the parameter Θ is of Lebesgue-measure zero. Therefore, a positive measure set allows you to query randomly and avoids getting coherent measurements.

Next two theorems identify the sufficient conditions for efficiently learning deterministic MDPs under polynomial realizability. Specifically, under online RL setting, we require the strong assumption that the set $\{\phi_h(s, a) | a \in \mathcal{A}\}$ is of positive measure for all $h \in [H]$ and all $s \in \mathcal{S}$, while under generative model setting, we only require the union set $\bigcup_{s \in \mathcal{S}} \{\phi_h(s, a) | a \in \mathcal{A}\}$ to be of positive measure for all $h \in [H]$. The algorithms for solving the both cases are summarized in Algorithms 4 and 5.

Assumption 4.6 (Polynomial Realizability). For all $h \in [H]$, $Q_h^*(s_h, a_h)$, viewed as the function of $\phi_h(s_h, a_h)$, lies in some admissible polynomial family \mathcal{F}_{ν_h} with dimension bounded by D .

Theorem 4.7. For the generative model setting, assume that the set $\{\phi_h(s, a) | s \in \mathcal{S}, a \in \mathcal{A}\}$ is of positive measure at any level h . Under the polynomial realizability, Algorithm 4 almost surely learns the optimal policy π^* with at most $N = 2DH$ samples.

Theorem 4.8. For the online RL setting, assume that $\{\phi_h(s, a) | a \in \mathcal{A}\}$ is of positive measure for every state s at every level h . Under polynomial realizability, within $T = 2DH$ episodes, Algorithm 5 learns the optimal policy π^* almost surely.

Algorithm 4 Dynamic programming under generative model settings

- 1: **for** $h = H, \dots, 1$ **do**
 - 2: Sample $2D$ points $\{\phi_h(s_h^{(i)}, a_h^{(i)})\}_{i=1}^{2D}$ according to $\mathbb{P}_{x \in \mathcal{N}(0, I_d)}(\cdot | x \in E_h)$ where $E_h = \{\phi_h(s, a) | s \in \mathcal{S}, a \in \mathcal{A}\}$.
 - 3: Query the generative model with state-action pair $(s_h^{(i)}, a_h^{(i)})$ at level h for $i = 1, \dots, 2D$, and observe the next state $\tilde{s}_h^{(i)}$ and reward $r_h^{(i)}$.
 - 4: Solve for Q_h^* with the $2D$ equations $Q_h^*(s_h^{(i)}, a_h^{(i)}) = r_h^{(i)} + V_{h+1}^*(\tilde{s}_h^{(i)})$.
 - 5: Set $\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$ and $V_h^*(s) = \max_a Q_h^*(s, a)$.
 - 6: **end for**
 - 7: **Output** π^*
-

Algorithm 5 Dynamic programming under online RL settings

- 1: **for** $h = H, \dots, 1$ **do**
 - 2: Fix any action sequence a_1, \dots, a_{h-1} .
 - 3: Play a_1, \dots, a_{h-1} for the first $h - 1$ levels and reach a state s_h . Sample $2D$ points $\{\phi_h(s_h, a_h^{(i)})\}_{i=1}^{2D}$ according to $\mathbb{P}_{x \in \mathcal{N}(0, I_d)}(\cdot | x \in E_h)$ where $E_h = \{\phi_h(s_h, a) | a \in \mathcal{A}\}$.
 - 4: Play $a_h^{(i)}$ at s_h for $i = 1, \dots, 2D$, and observe the next state $\tilde{s}_h^{(i)}$ and reward $r_h^{(i)}$.
 - 5: Solve for Q_h^* with the $2D$ equations $Q_h^*(s_h^{(i)}, a_h^{(i)}) = r_h^{(i)} + V_{h+1}^*(\tilde{s}_h^{(i)})$.
 - 6: Set $\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$ and $V_h^*(s) = \max_a Q_h^*(s, a)$.
 - 7: **end for**
 - 8: **Output** π^*
-

We remark that our Theorem 4.8 for learning MDPs under the online RL setting relies on a very strong assumption that allows the learner to explore randomly for any state. However, this assumption is necessary in some sense, as is suggested by our lower bound example in the next subsection.

4.1 Necessity of Generic Feature Maps in Online RL

In this section, we consider lower bounds for learning deterministic MDPs with polynomial realizable Q^* under online RL setting. Our goal is to show that in the online setting the generic assumption on

the feature maps $\phi_h(s, \cdot)$ is necessary. On the contrary, under the generative model setting one can efficiently learn the MDPs without such a strong assumption, since at every level h we can set the state arbitrarily.

MDP construction We briefly introduce the intuition of our construction. Consider a family of MDPs with only two states $\mathcal{S} = \{S_{\text{good}}, S_{\text{bad}}\}$. We set the feature map ϕ_h such that, for the good state S_{good} , it allows the learner to explore randomly, i.e., $\{\phi_h(S_{\text{good}}, a) \mid a \in \mathcal{A}\}$ is of positive measure.

However, for the bad state S_{bad} , all actions are mapped to some restricted set, which forbids random exploration, i.e., $\{\phi_h(S_{\text{bad}}, a) \mid a \in \mathcal{A}\}$ is measure zero. This is illustrated in Figure 1.

Specifically, at least $\Omega(d^p)$ actions are needed to identify the ground-truth polynomial of Q_h^* for S_{bad} , while $O(d)$ actions suffice for S_{good} .

The transition \mathbb{P}_h is constructed as $\mathbb{P}_h(s_{\text{bad}}|s, a) = 1$ for all $s \in \mathcal{S}, a \in \mathcal{A}$, which means it is impossible for the online scenarios to reach the good state for $h > 1$.

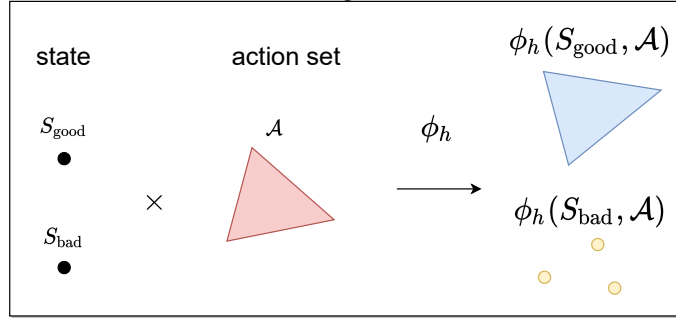


Figure 1: An illustration of the hard case for deterministic MDPs with polynomial realizable Q^* . The image of the feature map ϕ_h at S_{good} is of positive measure, while the image of ϕ_h at S_{bad} is not. This makes it difficult to learn under the online RL setting.

Theorem 4.9. *There exists a family of MDPs satisfying Assumption 4.6, such that the set $\{\phi_h(s, a) \mid s \in \mathcal{S}, a \in \mathcal{A}\}$ is of positive measure at any level h , but for all h there is some $s_{\text{bad}} \in \mathcal{S}$ such that $\{\phi_h(s_{\text{bad}}, a) \mid a \in \mathcal{A}\}$ is measure zero. Under the online RL setting, any algorithm needs to play at least $\Omega(d^p)$ episodes to identify the optimal policy. On the contrary, under the generative model setting, only $O(d)H$ samples are needed.*

5 Conclusions

In this paper, we consider neural network and polynomial function approximations in the simulator and online settings. To our knowledge, this is the first paper that shows sample-efficient reinforcement learning is possible with neural net function approximation. Our results substantially improve upon what can be achieved with existing results that primarily rely on embedding neural networks into linear function classes. The analysis reveals that for function approximations that allows for efficient ℓ_∞ recovery, such as two layer neural networks and admissible polynomial families, reinforcement learning can be reduced to parameter recovery problems, as well-studied in theories for deep learning, phase retrieval, and etc. Our method can also be potentially extended to handle three-layer and deeper neural networks, with advanced tools in [23, 24].

Our results for polynomial activation require deterministic transitions, since we cannot handle how noise propagates in solving polynomial equations. We leave to future work an in-depth study of the stability of roots of polynomial systems with noise, which is a fundamental mathematical problem and even unsolved for homogeneous polynomials. In particular, noisy tensor decomposition approaches combined with zeroth-order optimization may allow for stochastic transitions [36].

In the online RL setting, we can only show efficient learning under a very strong yet necessary assumption on the feature mapping. We leave to future work identifying more realistic and natural conditions which permit efficient learning in the online RL setting.

Finally, in future work, we hope to consider deep neural networks where parameter recovery or ℓ_∞ error is unattainable, and deep reinforcement learning with representation learning [84, 15].

Acknowledgment

JDL acknowledges support of the ARO under MURI Award W911NF-11-1-0303, the Sloan Research Fellowship, NSF CCF 2002272, and an ONR Young Investigator Award. QL is supported by NSF 2030859 and the Computing Research Association for the CIFellows Project. SK acknowledges funding from the NSF Award CCF-1703574 and the ONR award N00014-18-1-2247.

References

- [1] Yasin Abbasi-Yadkori, Nevena Lazic, Csaba Szepesvari, and Gellert Weisz. Exploration-enhanced POLITEX. *arXiv preprint arXiv:1908.10479*, 2019.
- [2] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [3] Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? *arXiv preprint arXiv:1905.10337*, 2019.
- [4] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning (ICML)*, pages 242–252, 2019.
- [5] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, pages 263–272, 2017.
- [6] Yu Bai and Jason D. Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. In *International Conference on Learning Representations*, 2020.
- [7] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, and Chris Hesse. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [8] Jacek Bochnak, Michel Coste, and Marie-Françoise Roy. *Real algebraic geometry*, volume 36. Springer Science & Business Media, 2013.
- [9] Minshuo Chen, Yu Bai, Jason D Lee, Tuo Zhao, Huan Wang, Caiming Xiong, and Richard Socher. Towards understanding hierarchical learning: Benefits of neural representations. *Neural Information Processing Systems (NeurIPS)*, 2020.
- [10] Alex Damian, Tengyu Ma, and Jason Lee. Label noise sgd provably prefers flat global minimizers. *arXiv preprint arXiv:2106.06530*, 2021.
- [11] Christoph Dann, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. On oracle-efficient PAC-RL with rich observations. In *Advances in Neural Information Processing Systems*, 2018.
- [12] Kefan Dong, Jian Peng, Yining Wang, and Yuan Zhou. \sqrt{n} -regret for learning in Markov decision processes with function approximation and low Bellman rank. In *Conference on Learning Theory*, pages 1554–1557. PMLR, 2020.
- [13] Kefan Dong, Jiaqi Yang, and Tengyu Ma. Provable model-based nonlinear bandit and reinforcement learning: Shelve optimism, embrace virtual curvature. *arXiv preprint arXiv:2102.04168*, 2021.
- [14] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning (ICML)*, pages 1675–1685. PMLR, 2019.
- [15] Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.

- [16] Simon S. Du, Sham M. Kakade, Jason D. Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in rl. *arXiv preprint arXiv:2103.10897*, 2021.
- [17] Simon S Du, Sham M Kakade, Jason D Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in RL. *arXiv preprint arXiv:2103.10897*, 2021.
- [18] Simon S. Du, Sham M. Kakade, Ruosong Wang, and Lin F. Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *International Conference on Learning Representations*, 2020.
- [19] Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *International Conference on Learning Representations*, 2020.
- [20] Simon S Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford. Provably efficient RL with rich observations via latent state decoding. In *International Conference on Machine Learning*, pages 1665–1674, 2019.
- [21] 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. *Advances in Neural Information Processing Systems*, 2020.
- [22] Cong Fang, Jason D Lee, Pengkun Yang, and Tong Zhang. Modeling from features: a mean-field framework for over-parameterized deep neural networks. *arXiv preprint arXiv:2007.01452*, 2020.
- [23] Massimo Fornasier, Timo Klock, and Michael Rauchensteiner. Robust and resource efficient identification of two hidden layer neural networks. *arXiv preprint arXiv:1907.00485*, 2019.
- [24] Massimo Fornasier, Jan Vybíral, and Ingrid Daubechies. Robust and resource efficient identification of shallow neural networks by fewest samples. *arXiv preprint arXiv:1804.01592*, 2019.
- [25] Dylan J. Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert E. Schapire. Practical contextual bandits with regression oracles. *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [26] Ruiqi Gao, Tianle Cai, Haochuan Li, Liwei Wang, Cho-Jui Hsieh, and Jason D Lee. Convergence of adversarial training in overparametrized networks. *Neural Information Processing Systems (NeurIPS)*, 2019.
- [27] Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- [28] Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *International Conference on Learning Representations (ICLR)*, 2018.
- [29] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. *International Conference on Learning Representations (ICLR)*, 2018.
- [30] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *arXiv preprint arXiv:1904.12191*, 2019.
- [31] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *The Annals of Statistics*, 49(2):1029–1054, 2021.
- [32] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [33] Steve Hanneke. Active learning for cost-sensitive classification. *Foundations and Trends® in Machine Learning*, 7(2-3):131–309, 2014.
- [34] Jeff Z. HaoChen, Colin Wei, Jason D. Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. *arXiv preprint arXiv:2006.08680*, 2020.

- [35] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [36] Baihe Huang, Kaixuan Huang, Sham M Kakade, Jason D Lee, Qi Lei, Runzhe Wang, and Jiaqi Yang. Optimal gradient-based algorithms for non-concave bandit optimization. *arXiv preprint arXiv:2107.04518*, 2021.
- [37] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems (NeurIPS)*, pages 8571–8580, 2018.
- [38] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- [39] Majid Janzamin, Hanie Sedghi, and Anankumar Anima. Beating the perils of nonconvexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- [40] Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1704–1713, 2017.
- [41] 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.
- [42] Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *arXiv preprint arXiv:2102.00875*, 2021.
- [43] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- [44] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143, 2020.
- [45] Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. PhD thesis, University of College London, 2003.
- [46] Akshay Krishnamurthy, Alekh Agarwal, Tzu-Kuo Huang, Hal Daume III, and John Langford. Active learning for cost-sensitive classification. *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [47] Volodymyr Kuleshov, Arun Chaganty, and Percy Liang. Tensor factorization via matrix factorization. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, page 507–516, 2015.
- [48] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [49] Gene Li, Pritish Kamath, Dylan J. Foster, and Nathan Srebro. Eluder dimension and generalized rank. *arXiv preprint arXiv:2104.06970*, 2021.
- [50] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [51] Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural proximal/trust region policy optimization attains globally optimal policy. *arXiv preprint arXiv:1906.10306*, 2019.
- [52] Dhruv Malik, Aldo Pacchiano, Vishwak Srinivasan, and Yuanzhi Li. Sample efficient reinforcement learning in continuous state spaces: A perspective beyond linearity. *arXiv preprint arXiv:2106.07814*, 2021.

- [53] James S. Milne. Algebraic geometry (v6.02), 2017. Available at www.jmilne.org/math/.
- [54] 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.
- [55] 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.
- [56] Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- [57] Edward Moroshko, Suriya Gunasekar, Blake Woodworth, Jason D Lee, Nathan Srebro, and Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy. *Neural Information Processing Systems (NeurIPS)*, 2020.
- [58] Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pages 4683–4692. PMLR, 2019.
- [59] Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems*, pages 3003–3011, 2013.
- [60] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning via sequential complexities. *Journal of Machine Learning Research*, 16(6):155–186, 2015.
- [61] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Sequential complexities and uniform martingale laws of large numbers. *Probability Theory and Related Fields*, 161(1-2):111–153, 2015.
- [62] Dan Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.
- [63] Daniel Russo. Worst-case regret bounds for exploration via randomized value functions. In *Advances in Neural Information Processing Systems*, pages 14433–14443, 2019.
- [64] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, 2013.
- [65] Igor R Shafarevich. *Basic Algebraic Geometry 1: Varieties in Projective Space*. Springer Science & Business Media, 2013.
- [66] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- [67] Aaron Sidford, Mengdi Wang, Xian Wu, Lin F Yang, and Yinyu Ye. Near-optimal time and sample complexities for solving discounted markov decision process with a generative model. *arXiv preprint arXiv:1806.01492*, 2018.
- [68] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, and Marc Lanctot. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484, 2016.
- [69] Alexander L Strehl, Lihong Li, Eric Wiewiora, John Langford, and Michael L Littman. PAC model-free reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 881–888. ACM, 2006.
- [70] 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.

- [71] Richard S Sutton and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [72] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2020.
- [73] Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Provably efficient reinforcement learning with general value function approximation. *arXiv preprint arXiv:2005.10804*, 2020.
- [74] Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Provably efficient reinforcement learning with general value function approximation. *Advances in Neural Information Processing Systems*, 2020.
- [75] Xiang Wang, Chenwei Wu, Jason D Lee, Tengyu Ma, and Rong Ge. Beyond lazy training for over-parameterized tensor decomposition. *Neural Information Processing Systems (NeurIPS)*, 2020.
- [76] Yang Wang and Zhiqiang Xu. Generalized phase retrieval: measurement number, matrix recovery and beyond. *Applied and Computational Harmonic Analysis*, 47(2):423–446, 2019.
- [77] Yining Wang, Ruosong Wang, Simon S. Du, and Akshay Krishnamurthy. Optimism in reinforcement learning with generalized linear function approximation. In *International Conference on Learning Representations*, 2021.
- [78] Yuanhao Wang, Ruosong Wang, and Sham M. Kakade. An exponential lower bound for linearly-realizable mdps with constant suboptimality gap. *arXiv preprint arXiv:2103.12690*, 2021.
- [79] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. In *Advances in Neural Information Processing Systems*, pages 9709–9721, 2019.
- [80] Gellert Weisz, Philip Amortila, and Csaba Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. *arXiv preprint arXiv:2010.01374*, 2020.
- [81] Zheng Wen and Benjamin Van Roy. Efficient exploration and value function generalization in deterministic systems. In *Advances in Neural Information Processing Systems*, pages 3021–3029, 2013.
- [82] Zheng Wen and Benjamin Van Roy. Efficient reinforcement learning in deterministic systems with value function generalization. *Math. Oper. Res.*, 42(3):762–782, 2017.
- [83] Blake Woodworth, Suriya Genesekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Kernel and deep regimes in overparametrized models. In *Conference on Learning Theory (COLT)*, 2019.
- [84] Jiaqi Yang, Wei Hu, Jason D Lee, and Simon S Du. Provable benefits of representation learning in linear bandits. *arXiv preprint arXiv:2010.06531*, 2020.
- [85] Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *International Conference on Machine Learning*, 2020.
- [86] Zhuoran Yang, Chi Jin, Zhaoran Wang, Mengdi Wang, and Michael Jordan. Provably efficient reinforcement learning with kernel and neural function approximations. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 13903–13916. Curran Associates, Inc., 2020.
- [87] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. *arXiv preprint arXiv:1904.00687*, 2019.
- [88] Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312, 2019.

- [89] Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent Bellman error. In *International Conference on Machine Learning*, 2020.
- [90] Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. *Proceedings of the Thirty-fourth International Conference on Machine Learning (ICML)*, 70, 2017.
- [91] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration. *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [92] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** Our abstract and introduction accurately reflect our main contribution.
 - (b) Did you describe the limitations of your work? **[Yes]** In the Preliminary section we clearly state the settings we considered and when our results apply.
 - (c) Did you discuss any potential negative societal impacts of your work? **[No]** Our work is theoretical and generally will not have negative social impacts.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** The assumptions are clearly included in every theorem or lemma.
 - (b) Did you include complete proofs of all theoretical results? **[Yes]** All the proofs can be found in the appendix
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[N/A]**
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[N/A]**
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[N/A]**
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[N/A]**
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? **[N/A]**
 - (b) Did you mention the license of the assets? **[N/A]**
 - (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[N/A]**
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]**
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

A Additional related work

Deterministic RL Deterministic system is often the starting case in the study of sample-efficient algorithms, where the issue of exploration and exploitation trade-off is more clearly revealed since both the transition kernel and reward function are deterministic. The seminal work [81] proposes a sample-efficient algorithm for Q-learning that works for a family of function classes. Recently, [21] studies the agnostic setting where the optimal Q-function can only be approximated by a function class with approximation error. The algorithm in [21] learns the optimal policy with the number of trajectories linear with the eluder dimension.

B Omitted Proofs in Section 3

B.1 Proof of Section 3.1

Theorem B.1 (Formal statement of Theorem 3.4). *Consider MDP \mathcal{M} where the transition is deterministic. Assume the function class in Definition 3.1 satisfies Assumption 2.1 and Assumption 2.2. For any $t \in (0, 1)$, if $d \geq \Omega(\log(B_W/\lambda))$ and $n \geq d \cdot \text{poly}(\kappa, k, \lambda, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ Algorithm 1 returns the optimal policy π^* .*

Proof. Use π_1^*, \dots, π_H^* to denote the global optimal policy. We prove that Algorithm 1 learns π_h^* from $h = H$ to $h = 1$.

At level H , the query obtains exact $Q_H^*(s, a)$. Therefore by Theorem B.15, $\widehat{Q}_H = Q_H^*$ and thus the optimal planning finds $\pi_H = \pi_H^*$. Suppose we have learned $\pi_{h+1}^*, \dots, \pi_H^*$ at level h . Due to deterministic transition, the query obtains exact $Q_h^*(s, a)$. Therefore by Theorem B.15, $\widehat{Q}_h = Q_h^*$ and thus the optimal planning finds $\pi_h = \pi_h^*$. Recursively applying this process to $h = 1$, we complete the proof. \square

B.2 Proof of Section 3.2

Theorem B.2 (Formal statement of Theorem 3.5). *Assume the function class in Definition 3.1 satisfies Assumption 2.1, Assumption 2.2 and is policy complete. For any $\epsilon > 0$ and $t \in (0, 1)$ such that $d \geq \Omega(\log(B_W B_\phi/\epsilon))$, if $n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ Algorithm 2 returns a policy π such that $V^* - V^\pi \leq \epsilon$.*

Proof. Use π_1^*, \dots, π_H^* to denote the global optimal policy. We prove for all $s \in \mathcal{S}$,

$$V_h^{\pi_h^*, \pi_{h+1}^*, \dots, \pi_H^*}(s) - V_h^{\pi_h, \pi_{h+1}, \dots, \pi_H}(s) \leq \frac{(H - h + 1)\epsilon}{H}.$$

At level H , let $e_H(s_H^i, a_H^i) = r_H(s_H^i, a_H^i) - Q_H^*(s_H^i, a_H^i)$, then $e_H(s_H^i, a_H^i) = 0$. From Theorem B.13, we have $\widehat{Q}_H(s, a) := v_H^\top \sigma(W_H \phi(s, a))$ satisfies $|\widehat{Q}_H(s, a) - Q_H^*(s, a)| \leq \frac{\epsilon}{2H}$ for all $s \in \mathcal{S}, a \in \mathcal{A}$. Therefore for all $s \in \mathcal{S}$,

$$\begin{aligned} V_H^*(s) - V_H^{\pi_H}(s) &= \mathbb{E}_{a \sim \pi_H^*} [Q_H^*(s, a)] - \mathbb{E}_{a \sim \pi_H^*} [\widehat{Q}_H(s, a)] \\ &\quad + \mathbb{E}_{a \sim \pi_H^*} [\widehat{Q}_H(s, a)] - \mathbb{E}_{a \sim \pi_H} [\widehat{Q}_H(s, a)] \\ &\quad + \mathbb{E}_{a \sim \pi_H} [\widehat{Q}_H(s, a)] - \mathbb{E}_{a \sim \pi_H} [Q_H^*(s, a)] \\ &\leq \frac{\epsilon}{H} \end{aligned}$$

where in the second step we used $\mathbb{E}_{a \sim \pi_H^*} [\widehat{Q}_H(s, a)] \leq \mathbb{E}_{a \sim \pi_H} [\widehat{Q}_H(s, a)]$ by optimality of π_H and $|\widehat{Q}_H(s, a) - Q_H^*(s, a)| \leq \frac{\epsilon}{2H}$.

Suppose we have learned policies π_{h+1}, \dots, π_H , we use $\tilde{\pi}_h$ to denote the optimal policy of $Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a)$. Let

$$e_h(s_h^i, a_h^i) = \widehat{Q}_h^i - Q_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i)$$

then $e_h(s_h^i, a_h^i)$ is zero mean H^2 sub-Gaussian (notice that \widehat{Q}_h^i is unbiased estimate of $Q_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i)$, and $\widehat{Q}_h^i \leq O(H)$). From Theorem B.13, we have $\widehat{Q}_h(s, a) = v_h^\top \sigma(W_h \phi(s, a))$ satisfies $|\widehat{Q}_h(s, a) - Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a)| \leq \frac{\epsilon}{2H}$ for all $s \in \mathcal{S}, a \in \mathcal{A}$. Therefore for all $s \in \mathcal{S}$,

$$\begin{aligned} & V_h^{\widetilde{\pi}_h, \pi_{h+1}, \dots, \pi_H}(s) - V_h^{\pi_h, \pi_{h+1}, \dots, \pi_H}(s) \\ &= \mathbb{E}_{a \sim \widetilde{\pi}_h} [Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a)] - \mathbb{E}_{a \sim \widetilde{\pi}_h} [\widehat{Q}_h(s, a)] \\ &\quad + \mathbb{E}_{a \sim \widetilde{\pi}_h} [\widehat{Q}_h(s, a)] - \mathbb{E}_{a \sim \pi_h} [\widehat{Q}_h(s, a)] \\ &\quad + \mathbb{E}_{a \sim \pi_h} [\widehat{Q}_h(s, a)] - \mathbb{E}_{a \sim \pi_h} [Q_h^{\pi_h, \pi_{h+1}, \dots, \pi_H}(s, a)] \\ &\leq \frac{\epsilon}{H} \end{aligned}$$

where in the second step we used $\mathbb{E}_{a \sim \widetilde{\pi}_h} [\widehat{Q}_h(s, a)] \leq \mathbb{E}_{a \sim \pi_h} [\widehat{Q}_h(s, a)]$ by optimality of π_h and $|\widehat{Q}_h(s, a) - Q_h^{\pi_{h+1}, \dots, \pi_H}(s, a)| \leq \frac{\epsilon}{2H}$.

It thus follows that

$$\begin{aligned} V_h^{\pi_h^*, \pi_{h+1}^*, \dots, \pi_H^*}(s) - V_h^{\pi_h, \pi_{h+1}, \dots, \pi_H}(s) &= V_h^{\pi_h^*, \pi_{h+1}^*, \dots, \pi_H^*}(s) - V_h^{\pi_h^*, \pi_{h+1}, \dots, \pi_H}(s) \\ &\quad + V_h^{\pi_h^*, \pi_{h+1}, \dots, \pi_H}(s) - V_h^{\widetilde{\pi}_h, \pi_{h+1}, \dots, \pi_H}(s) \\ &\quad + V_h^{\widetilde{\pi}_h, \pi_{h+1}, \dots, \pi_H}(s) - V_h^{\pi_h, \pi_{h+1}, \dots, \pi_H}(s) \\ &\leq V_h^{\pi_h^*, \pi_{h+1}^*, \dots, \pi_H^*}(s) - V_h^{\pi_h^*, \pi_{h+1}, \dots, \pi_H}(s) + \frac{\epsilon}{H} \\ &\leq \dots \\ &\leq \frac{(H-h+1)\epsilon}{H}. \end{aligned}$$

where in the second step we use $V_h^{\pi_h^*, \pi_{h+1}, \dots, \pi_H}(s) \leq V_h^{\widetilde{\pi}_h, \pi_{h+1}, \dots, \pi_H}(s)$ from optimality of $\widetilde{\pi}_h$. Repeating this argument to $h=1$ completes the proof \square

B.3 Proof of Section 3.3

Theorem B.3 (Formal statement of Theorem 3.6). *Assume the function class in Definition 3.1 satisfies Assumption 2.1, Assumption 2.2, and is Bellman complete. For any $\epsilon > 0$ and $t \in (0, 1)$ such that $d \geq \Omega(\log(B_W B_\phi / \epsilon))$, if $n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1-t$ Algorithm 3 returns a policy π such that $V^* - V^\pi \leq \epsilon$.*

Proof. Use π_1^*, \dots, π_H^* to denote the global optimal policy. We prove

$$|\widehat{Q}_h(s, a) - Q_h^*(s, a)| \leq \frac{(H-h+1)\epsilon}{H} \quad (1)$$

for all $s \in \mathcal{S}, a \in \mathcal{A}$.

At level H , let

$$e_H(s_H^i, a_H^i) = r_H(s_H^i, a_H^i) - Q_H^*(s_H^i, a_H^i)$$

then $e_H(s_H^i, a_H^i) = 0$. From Theorem B.13, we have $\widehat{Q}_H(s, a) := v_H^\top \sigma(W_H \phi(s, a))$ satisfies $|\widehat{Q}_H(s, a) - Q_H^*(s, a)| \leq \frac{\epsilon}{H}$ for all $s \in \mathcal{S}, a \in \mathcal{A}$.

Suppose we have learned $\widehat{Q}_{h+1}(s, a)$ with $|\widehat{Q}_{h+1}(s, a) - Q_{h+1}^*(s, a)| \leq \frac{(H-h)\epsilon}{H}$. At level h , let

$$e_h(s_h^i, a_h^i) = r_h(s_h^i, a_h^i) + \widehat{V}_{h+1}(s_{h+1}^i) - \mathcal{T}_h(\widehat{Q}_{h+1})(s_h^i, a_h^i)$$

then $e_h(s_h^i, a_h^i)$ is zero mean H^2 sub-Gaussian (notice that $r_h(s_h^i, a_h^i) + \widehat{V}_{h+1}(s_{h+1}^i)$ is unbiased estimate of $\mathcal{T}_h(\widehat{Q}_{h+1})(s_h^i, a_h^i)$, and $r_h(s_h^i, a_h^i) + \widehat{V}_{h+1}(s_{h+1}^i) \leq O(H)$). From Theorem B.13, we have $\widehat{Q}_h(s, a) := v_h^\top \sigma(W_h \phi(s, a))$ satisfies $|\widehat{Q}_h(s, a) - \mathcal{T}_h(\widehat{Q}_{h+1})(s_h^i, a_h^i)| \leq \frac{\epsilon}{H}$ for all $s \in \mathcal{S}, a \in \mathcal{A}$.

A. Therefore

$$\begin{aligned}
|\widehat{Q}_h(s, a) - Q_h^*(s, a)| &\leq |\widehat{Q}_h(s, a) - \mathcal{T}_h(\widehat{Q}_{h+1})(s, a)| + |\mathcal{T}_h(\widehat{Q}_{h+1})(s, a) - Q_h^*(s, a)| \\
&\leq \frac{\epsilon}{H} + \max_{s \in \mathcal{S}, a \in \mathcal{A}} |\widehat{Q}_{h+1}(s, a) - Q_{h+1}^*(s, a)| \\
&\leq \frac{(H - h + 1)\epsilon}{H}
\end{aligned}$$

holds for all $s \in \mathcal{S}, a \in \mathcal{A}$.

It thus follows that for all $s_1 \in \mathcal{S}$,

$$\begin{aligned}
V_h^{\pi_1^*, \dots, \pi_H^*}(s_1) - V_h^{\pi_1, \dots, \pi_H}(s_1) &= \mathbb{E}_{a \sim \pi_1^*}[Q_1^*(s_1, a)] - \mathbb{E}_{a \sim \pi_1}[Q_1^{\pi_2, \dots, \pi_H}(s_1, a)] \\
&\leq \mathbb{E}_{a \sim \pi_1^*}[\widehat{Q}_1(s_1, a)] - \mathbb{E}_{a \sim \pi_1}[Q_1^{\pi_2, \dots, \pi_H}(s_1, a)] + \epsilon \\
&\leq \mathbb{E}_{a \sim \pi_1}[\widehat{Q}_1(s_1, a) - Q_1^{\pi_2, \dots, \pi_H}(s_1, a)] + \epsilon \\
&\leq \mathbb{E}_{a \sim \pi_1}[Q_1^*(s_1, a) - Q_1^{\pi_2, \dots, \pi_H}(s_1, a)] + 2\epsilon \\
&\leq \mathbb{E}_{a \sim \pi_1} \mathbb{E}_{s_2 \sim \mathbb{P}(\cdot | s_1, a)}[V_2^{\pi_2^*, \dots, \pi_H^*}(s_2) - V_2^{\pi_2, \dots, \pi_H}(s_2)] + 2\epsilon \\
&\leq \dots \\
&\leq 2H\epsilon
\end{aligned}$$

where the first step comes from definition of value function; the second step comes from Eq (1); the third step comes from optimality of π_1 ; the fourth step comes from Eq (1); the fifth step comes from Bellman equation. The proof is complete by rescaling $\epsilon \leftarrow \epsilon/H$. \square

B.4 Proof of Section 3.4

With gap condition, either Algorithm 2 or Algorithm 3 will work as long as we select $\epsilon \approx \rho$. The following displays an adaption from Algorithm 2.

Algorithm 6 Learning realizable Q^* with optimality gap

```

1: for  $h = H, \dots, 1$  do
2:   Sample  $x_h^i, i \in [n]$  from standard Gaussian  $\mathcal{N}(0, I_d)$ 
3:   for  $i \in [n]$  do
4:     if  $\|x_h^i\| \leq \delta_\phi$  then
5:       Find  $(s_h^i, a_h^i) \in \phi^{-1}(x_h^i)$  and locate the state  $s_h^i$  in the generative model
6:       Pull action  $a_h^i$  and use  $\pi_{h+1}, \dots, \pi_H$  as the roll-out to collect rewards  $r_h^{(i)}, \dots, r_H^{(i)}$ 
7:       Construct unbiased estimation of  $Q_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i)$ 

$$\widehat{Q}_h^i \leftarrow r_h^{(i)} + \dots + r_H^{(i)}$$

8:     else
9:       Let  $\widehat{Q}_h^i \leftarrow 0$ .
10:    end if
11:  end for
12:  Compute  $(v_h, W_h) \leftarrow \text{NEURALNETNOISYRECOVERY}(\{(x_h^i, \widehat{Q}_h^i) : i \in [n]\})$ 
13:  Set  $\widehat{Q}_h(s, a) \leftarrow v_h^\top \sigma(W_h \phi(s, a))$ 
14:  Let  $\pi_h(s) \leftarrow \arg \max_{a \in \mathcal{S}} \widehat{Q}_h(s, a)$ 
15: end for
16: Return  $\pi_1, \dots, \pi_H$ 

```

Theorem B.4 (Formal statement of Theorem 3.8). *Assume the function class in Definition 3.1 satisfies Assumption 2.1 and Assumption 2.2. Suppose $\rho > 0$ and $d \geq \Omega(\log(B_W B_\phi / \rho))$, for any $t \in (0, 1)$, if $n = \frac{d}{\rho^2} \cdot \text{poly}(\kappa, k, B_W, B_\phi, H, \log(d/t))$, then with probability at least $1 - t$ Algorithm 6 returns the optimal policy π^* .*

Proof. Use π_1^*, \dots, π_H^* to denote the global optimal policy. Similar to Theorem B.1, we prove that Algorithm 6 learns π_h^* from $h = H$ to $h = 1$.

At level H , the algorithm uses $n = \frac{d}{\rho^2} \cdot \text{poly}(\kappa, k, \log d, B_W, B_\phi, H, \log(1/t))$ trajectories to obtain \widehat{Q}_H such that $|\widehat{Q}_H(s, a) - Q^*(s, a)| \leq \rho/4$ by Theorem B.13. Therefore

$$\begin{aligned} & V_H^*(s) - Q_H^*(s, \pi_H(s)) \\ & \leq Q_H^*(s, \pi_H^*(s)) - Q_H^*(s, \pi_H(s)) \\ & \leq Q_H^*(s, \pi_H^*(s)) - \widehat{Q}_H(s, \pi_H^*(s)) + \widehat{Q}_H(s, \pi_H^*(s)) - \widehat{Q}_H(s, \pi_H(s)) \\ & \quad + \widehat{Q}_H(s, \pi_H(s)) - Q_H^*(s, \pi_H(s)) \\ & \leq \rho/2 \end{aligned}$$

where the third inequality uses the optimality of $\pi_H(s)$ under \widehat{Q}_H . Thus Definition 3.7 gives $\pi_H(s) = \pi_H^*(s)$. Suppose we have learned $\pi_{h+1}^*, \dots, \pi_H^*$ at level h . We apply the same argument to derive $\pi_h = \pi_h^*$. Recursively applying this process to $h = 1$, we complete the proof. \square

B.5 Neural network recovery

This section considers recovering neural network $\langle v, \sigma(Wx) \rangle$ from the following two models, where $B = \Omega(d \cdot \text{poly} \log(d))$.

- Noisy samples from

$$x \sim \mathcal{N}(0, I_d), \quad y = (\langle v, \sigma(Wx) \rangle + \xi) \cdot \mathbb{1}(\|x\| \leq B) \quad (2)$$

where ξ is ϑ sub-Gaussian noise.

- Noiseless samples from

$$x \sim \mathcal{N}(0, I_d), \quad y = (\langle v, \sigma(Wx) \rangle) \cdot \mathbb{1}(\|x\| \leq B) \quad (3)$$

Recovering neural network has received comprehensive study in deep learning theory [39, 90, 27]. The analysis in this section is mainly based on the method of moments in [90]. However, notice that the above learning tasks are different from those considered in [90], due to the presence of noise and the truncated signals. Therefore, additional considerations must be made in the analysis.

We consider more general homogeneous activation functions, specified by the assumptions that follow. Since the activation function is homogeneous, we assume $v_i \in \{\pm 1\}$ in the following without loss of generality.

Assumption B.5 (Property 3.1 of [90]). Assume $\sigma'(x)$ is nonnegative and homogeneously bounded, i.e. $0 \leq \sigma'(x) \leq L_1|x|^p$ for some constants $L_1 > 0$ and $p \geq 0$.

Definition B.6 (Part of property 3.2 of [90]). Define $\rho(z) := \min\{\beta_0(z) - \alpha_0^2(z) - \alpha_1^2(z), \beta_2(z) - \alpha_1^2(z) - \alpha_2^2(z), \alpha_0(z)\alpha_2(z) - \alpha_1^2(z)\}$, where $\alpha_q(z) := \mathbb{E}_{x \sim \mathcal{N}(0,1)}[\sigma'(zx)x^q]$, $q \in \{0, 1, 2\}$, and $\beta_q(z) := \mathbb{E}_{x \sim \mathcal{N}(0,1)}[(\sigma')^2(zx)x^q]$ for $q \in \{0, 2\}$.

Assumption B.7 (Part of property 3.2 of [90]). The first derivative $\sigma'(z)$ satisfies that, for all $z > 0$, we have $\rho(z) > 0$.

Assumption B.8 (Property 3.3 of [90]). The second derivative $\sigma''(x)$ is either **(a)** globally bounded or **(b)** $\sigma''(x) = 0$ except for finite points.

Notice that ReLU, squared ReLU, leaky ReLU, and polynomial activation function functions all satisfies the above assumption. We make the following assumption on the dimension of feature vectors, which corresponds to how features can extract information about neural networks from noisy samples. The dimension only has to be greater than a logarithmic term in $1/\epsilon$ and the norm of parameters.

Assumption B.9 (Rich feature). Assume $d \geq \Omega(\log(B_W/\epsilon))$.

First we introduce a notation from [90].

Definition B.10. Define outer product $\tilde{\otimes}$ as follows. For a vector $v \in \mathbb{R}^d$ and an identity matrix $I \in \mathbb{R}^{d \times d}$,

$$v \tilde{\otimes} I = \sum_{j=1}^d [v \otimes e_j \otimes e_j + e_j \otimes v \otimes e_j + e_j \otimes e_j \otimes v].$$

For a symmetric rank- r matrix $M = \sum_{i=1}^r s_i v_i v_i^\top$ and an identity matrix $I \in \mathbb{R}^{d \times d}$,

$$M \tilde{\otimes} I = \sum_{i=1}^r s_i \sum_{j=1}^d \sum_{l=1}^6 A_{l,i,j}$$

where $A_{1,i,j} = v_i \otimes v_i \otimes e_j \otimes e_j$, $A_{2,i,j} = v_i \otimes e_j \otimes v_i \otimes e_j$, $A_{3,i,j} = e_j \otimes v_i \otimes v_i \otimes e_j$, $A_{4,i,j} = v_i \otimes e_j \otimes e_j \otimes v_i$, $A_{5,i,j} = e_j \otimes v_i \otimes e_j \otimes v_i$, $A_{6,i,j} = e_j \otimes e_j \otimes v_i \otimes v_i$.

Now we define some moments.

Definition B.11. Define $M_1, M_2, M_3, M_4, m_{1,i}, m_{2,i}, m_{3,i}, m_{4,i}$ as follows:

$$\begin{aligned} M_1 &:= \mathbb{E}[y \cdot x] \\ M_2 &:= \mathbb{E}[y \cdot (x \otimes x - I)] \\ M_3 &:= \mathbb{E}[y \cdot (x^{\otimes 3} - x \tilde{\otimes} I)] \\ M_4 &:= \mathbb{E}[y \cdot (x^{\otimes 4} - (x \otimes x) \tilde{\otimes} I + I \tilde{\otimes} I)] \\ \gamma_j(x) &:= \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma(x \cdot z) z^j], \forall j \in \{0, 1, 2, 3, 4\} \\ m_{1,i} &:= \gamma_1(\|w_i\|) \\ m_{2,i} &:= \gamma_2(\|w_i\|) - \gamma_0(\|w_i\|) \\ m_{3,i} &:= \gamma_3(\|w_i\|) - 3\gamma_1(\|w_i\|) \\ m_{4,i} &:= \gamma_4(\|w_i\|) + 3\gamma_0(\|w_i\|) - 6\gamma_2(\|w_i\|) \end{aligned}$$

The above expectations are all with respect to $x \sim \mathcal{N}(0, I_d)$ and $y = \langle v, \sigma(Wx) \rangle$.

Assumption B.12 (Assumption 5.3 of [90]). Assume the activation function satisfies the followings:

- If $M_i \neq 0$, then $m_{j,i} \neq 0$ for all $i \in [k]$.
- At least one of M_3 and M_4 is not zero.
- If $M_1 = M_3 = 0$, then $\sigma(z)$ is an even function.
- If $M_2 = M_4 = 0$, then $\sigma(z)$ is an odd function.

Now we state the theoretical result that recovers neural networks from noisy data.

Theorem B.13 (Neural network recovery from noisy data). *Let the activation function σ satisfies Assumption B.5 and Assumption B.12. Let κ be the condition number of W . Given n samples from Eq (2). For any $t \in (0, 1)$ and $\epsilon \in (0, 1)$ such that Assumption B.9 holds, if*

$$n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, \vartheta, \log(d/t))$$

then there exists an algorithm that takes $\tilde{O}(nkd)$ time and returns a matrix $\widehat{W} \in \mathbb{R}^{k \times d}$ and a vector $\widehat{v} \in \{\pm 1\}^k$ such that with probability at least $1 - t$,

$$\|\widehat{W} - W\|_F \leq \epsilon \cdot \text{poly}(k, \kappa) \cdot \|W\|_F, \text{ and } \widehat{v} = v.$$

The algorithm and proof are shown in Appendix B.5.1. By Assumption B.5, the following corollary is therefore straightforward.

Corollary B.14. *In the same setting as Theorem B.13. For any $t \in (0, 1)$ and suppose $\|W\|_F \leq B_W$ and Assumption B.9 holds. Given n samples from Eq (2). If*

$$n \geq \epsilon^{-2} \cdot d \cdot \text{poly}(\kappa, k, \log(d/t), B_W, B_\phi, \vartheta)$$

then there exists an algorithm that takes $\tilde{O}(nkd)$ time and outputs a matrix $\widehat{W} \in \mathbb{R}^{k \times d}$ and a vector $\widehat{v} \in \{\pm 1\}^k$ such that with probability at least $1 - t$, for all $\|x\|_2 \leq B_\phi$

$$|\langle \widehat{v}, \sigma(\widehat{W}x) \rangle - \langle v, \sigma(Wx) \rangle| \leq \epsilon.$$

In particular, when $B_\phi = O(d \cdot \text{poly} \log d)$ the following sample complexity suffices

$$n \geq \epsilon^{-2} \cdot d^{O(1+p)} \cdot \text{poly}(\kappa, k, \log(d/t), B_W, \vartheta).$$

Now we state the theoretical result that precisely recovers neural networks from noiseless data. The proof and method are shown in Appendix B.5.2.

Theorem B.15 (Exact neural network recovery from noiseless data). *Let the activation function satisfies Assumption B.5 and Assumption B.12, Assumption B.7 and Assumption B.8(b). Given n samples from Eq (3). For any $t \in (0, 1)$, suppose $d \geq \Omega(\log(B_W/\lambda))$ and*

$$n \geq d \cdot \text{poly}(\kappa, k, \lambda, \log(d/t)),$$

then there exists an algorithm that output exact W and v with probability at least $1 - t$.

B.5.1 Recover neural networks from noisy data

In this section we prove Theorem B.13. Denote $W = [w_1, \dots, w_k]^\top$ where $w_i \in \mathbb{R}^d$ and $\bar{w}_i = w_i / \|w_i\|_2$.

Definition B.16. Given a vector $\alpha \in \mathbb{R}^d$. Define $P_2 := M_{j_2}(I, I, \alpha, \dots, \alpha)$ where $j_2 = \min\{j \geq 2 : M_j \neq 0\}$ and $P_3 := M_{j_3}(I, I, I, \alpha, \dots, \alpha)$ where $j_3 = \min\{j \geq 3 : M_j \neq 0\}$.

The method of moments is presented in Algorithm 7. Here we sketch its ideas, and refer readers to [90] for thorough explanations. There are three main steps. In the first step, it computes the span of the rows of W . By power method, Line 7 finds the top- k eigenvalues of $CI + \hat{P}_2$ and $CI - \hat{P}_2$. It then picks the largest k eigenvalues from $CI + \hat{P}_2$ and $CI - \hat{P}_2$, by invoking TOPK in Line 15. Finally it orthogonalizes the corresponding eigenvectors in Line 19 and finds an orthogonal matrix V in the subspace spanned by $\{\bar{w}_1, \dots, \bar{w}_k\}$.

In the second step, the algorithm forms third order tensor $R_3 = P_s(V, V, V) \in \mathbb{R}^{k \times k \times k}$ and use the robust tensor decomposition method in [47] to find \hat{u} that approximates $s_i V^\top \bar{w}_i$ with unknown signs s_i . In the third step, the algorithm determines s, v and $w_i, i \in [k]$. Since the activation function is homogeneous, we assume $v_i \in \{\pm 1\}$ and $m_{j,i} = c_j \|w_i\|^{p+1}$ for universal constants c_j without loss of generality. For illustration, we define Q_1 and Q_2 as follows.

$$Q_1 = M_{l_1}(I, \underbrace{\alpha, \dots, \alpha}_{(l_1-1) \text{ } \alpha\text{'s}}) = \sum_{i=1}^k v_i c_{l_1} \|w_i\|^{p+1} (\alpha^\top \bar{w}_i)^{l_1-1} \bar{w}_i, \quad (4)$$

$$Q_2 = M_{l_2}(V, V, \underbrace{\alpha, \dots, \alpha}_{(l_2-2) \text{ } \alpha\text{'s}}) = \sum_{i=1}^k v_i c_{l_2} \|w_i\|^{p+1} (\alpha^\top \bar{w}_i)^{l_2-2} (V^\top \bar{w}_i) (V^\top \bar{w}_i)^\top, \quad (5)$$

where $l_1 \geq 1$ such that $M_{l_1} \neq 0$ and $l_2 \geq 2$ such that $M_{l_2} \neq 0$ are specified later. Then the solutions of the following linear systems

$$z^* = \arg \min_{z \in \mathbb{R}^k} \left\| \sum_{i=1}^k z_i s_i \bar{w}_i - Q_1 \right\|, \quad r^* = \arg \min_{r \in \mathbb{R}^k} \left\| \sum_{i=1}^k r_i V^\top \bar{w}_i (V^\top \bar{w}_i)^\top - Q_2 \right\|_F. \quad (6)$$

are the followings

$$z_i^* = v_i s_i^{l_1} c_{l_1} \|w_i\|^{p+1} (\alpha^\top s_i \bar{w}_i)^{l_1-1}, \quad r_i = v_i s_i^{l_2} c_{l_2} \|w_i\|^{p+1} (\alpha^\top s_i \bar{w}_i)^{l_2-2}.$$

When c_{l_1} and c_{l_2} do not have the same sign, we can recover v_i and s_i by $v_i = \text{sign}(r_i^* c_{l_2}), s_i = \text{sign}(v_i z_i^* c_{l_1})$, and recover w_i by

$$w_i = \left(\left| \frac{z_i^*}{c_{l_1} (\alpha^\top s_i \bar{w}_i)^{l_1-1}} \right| \right)^{1/(p+1)} \bar{w}_i.$$

In Algorithm 7, we use $V \hat{u}_i$ to approximate $s_i \bar{w}_i$, and use moment estimators \hat{Q}_1 and \hat{Q}_2 to approximate Q_1 and Q_2 . Then the solutions \hat{z}, \hat{r} to the optimization problems in Line 29 should approximate z^* and r^* , due to robustness for solving linear systems. As such, the outputs \tilde{v}, \tilde{W} approximately recover the true model parameters.

Since Algorithm 7 carries out the same computation as [90], the computational complexity is the same. The difference of sample complexity comes from the noise ξ in the model and the truncation of standard Gaussian. The proof entails bounding the error in estimating P_2 in Line 4, R_3 in Line 20 and Q_1, Q_2 in Line 28. In the following, unless further specified, the expectations are all with respect to $x \sim \mathcal{N}(0, I_d)$ and $y \sim (\langle v, \sigma(Wx) \rangle + \xi) \cdot \mathbb{1}(\|x\| \leq B)$.

Algorithm 7 Using method of moments to recover neural network parameters

```

1: procedure NEURALNETNOISYRECOVERY( $S = \{(x_i, y_i) : i \in [n]\}$ )
2:   Choose  $\alpha$  to be a random unit vector
3:   Partition  $S$  into  $S_1, S_2, S_3, S_4$  of equal size
4:    $\widehat{P}_2 \leftarrow \mathbb{E}_{S_1}[P_2]$ ,  $C \leftarrow 3\|P_2\|$ ,  $T \leftarrow \log(1/\epsilon)$ 
5:   Choose  $\widehat{V}_1^{(0)}, \widehat{V}_2^{(0)} \in \mathbb{R}^{d \times k}$  to be random matrices ▷ Estimate subspace  $V$ 
6:   for  $t = 1, \dots, T$  do
7:      $\widehat{V}_1^{(t)} \leftarrow \text{QR}(C\widehat{V}_1^{(t-1)} + \widehat{P}_2\widehat{V}_1^{(t-1)})$ ,  $\widehat{V}_2^{(t)} \leftarrow \text{QR}(C\widehat{V}_2^{(t-1)} - \widehat{P}_2\widehat{V}_2^{(t-1)})$ 
8:   end for
9:   for  $j = 1, 2$  do
10:     $\widehat{V}_1^{(T)} \leftarrow [\widehat{V}_{j,1}, \dots, \widehat{V}_{j,k}]$ 
11:    for  $i \in [k]$  do
12:       $\lambda_{j,i} \leftarrow |\widehat{V}_{j,i} \widehat{P}_2 \widehat{V}_{j,i}|$ 
13:    end for
14:  end for
15:   $\pi_1, \pi_2, k_1, k_2 \leftarrow \text{TOPK}(\lambda, k)$ 
16:  for  $j = 1, 2$  do
17:     $V_j \leftarrow [\widehat{V}_{j,\pi_j(1)}, \dots, \widehat{V}_{j,\pi_j(k_j)}]$ 
18:  end for
19:   $\widetilde{V}_2 \leftarrow \text{QR}((I - V_1 V_1^\top) V_2)$ ,  $V \leftarrow [V_1, \widetilde{V}_2]$ 
20:   $\widehat{R}_3 \leftarrow \mathbb{E}_{S_2}[P_3(V, V, V)]$ ,  $\{\widehat{u}_i\}_{i \in [k]} \leftarrow \text{TensorDecomposition}(\widehat{R}_3)$  ▷ Learn  $s_i V^\top \bar{w}_i$ 
21:  if  $M_1 = M_3 = 0$  then
22:     $l_1, l_2 = \min\{j \in \{2, 4\} : M_j \neq 0\}$ 
23:  else if  $M_2 = M_4 = 0$  then
24:     $l_1 \leftarrow \min\{j \in \{1, 3\} : M_j \neq 0\}$ ,  $l_2 \leftarrow 3$ 
25:  else
26:     $l_1 \leftarrow \min\{j \in \{1, 3\} : M_j \neq 0\}$ ,  $l_2 = \min\{j \in \{2, 4\} : M_j \neq 0\}$ 
27:  end if
28:   $\widehat{Q}_1 \leftarrow \mathbb{E}_{S_3}[Q_1]$ ,  $\widehat{Q}_2 \leftarrow \mathbb{E}_{S_4}[Q_2]$ 
29:   $\widehat{z} \leftarrow \arg \min_z \|\sum_{i=1}^k z_i V \widehat{u}_i - \widehat{Q}_1\|$ ,  $\widehat{r} \leftarrow \arg \min_r \|\sum_{i=1}^k r_i \widehat{u}_i \widehat{u}_i - \widehat{Q}_2\|_F$ 
30:  for  $i = 1, \dots, k$  do ▷ Learn parameters  $v, W$ 
31:     $\widehat{v}_i \leftarrow \text{sign}(\widehat{r}_i c_{l_2})$ ,  $\widehat{s}_i \leftarrow \text{sign}(\widehat{v}_i \widehat{z}_i c_{l_1})$ 
32:     $\widehat{w}_i \leftarrow \widehat{s}_i (|\frac{\widehat{z}_i}{c_{l_1} (\alpha^\top V \widehat{u}_i)^{l_1 - 1}}|)^{1/(p+1)} V \widehat{u}_i$ 
33:  end for
34:   $\widehat{W} \leftarrow [\widehat{w}_1, \dots, \widehat{w}_k]$ 
35:  Return  $(\widehat{v}, \widehat{W})$ 
36: end procedure

```

Lemma B.17. Let \widehat{P}_2 be computed in Line 4 of Algorithm 7 and P_2 defined in Definition B.16. Suppose $m_0 = \min_{i \in [k]} \{|m_{j_2, i}|^2 (\bar{w}_i^\top \alpha)^{2(j_2 - 2)}\}$ and

$$|S| \gtrsim d \cdot \text{poly}(\kappa, \vartheta, \log(d/t)) / (\epsilon^2 m_0)$$

then with probability at least $1 - t$,

$$\|P_2 - \widehat{P}_2\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{j_2, i} (\bar{w}_i^\top \alpha)^{j_2 - 2}| + \epsilon.$$

Proof. It suffices to bound $\|M_2 - \widehat{M}_2\|$, $\|M_3(I, I, \alpha) - \widehat{M}_3(I, I, \alpha)\|$ and $\|M_4(I, I, \alpha, \alpha) - \widehat{M}_4(I, I, \alpha, \alpha)\|$. The main strategy is to bound all relevant moment terms and to invoke Claim E.6. Specifically, we show that with probability at least $1 - t/4$,

$$\|M_2 - \widehat{M}_2\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{2, i}| + \epsilon. \tag{7}$$

$$\|M_3(I, I, \alpha) - \widehat{M}_3(I, I, \alpha)\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{3,i}(\bar{w}_i^\top \alpha)| + \epsilon. \quad (8)$$

$$\|M_4(I, I, \alpha, \alpha) - \widehat{M}_4(I, I, \alpha, \alpha)\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{4,i}|(\bar{w}_i^\top \alpha)^2 + \epsilon. \quad (9)$$

Recall that for sample $(x_j, y_j) \in S$, $y_j = \sum_{i=1}^k v_i \sigma(w_i^\top x_j) + \xi_j$ where ξ_j is independent of x_j . Consider each component $i \in [k]$. Define $C_i(x_j), B_i(x_j) \in \mathbb{R}^{d \times d}$ as follows:

$$\begin{aligned} B_i(x_j) &= (\sigma(w_i^\top x_j) + \xi_j) \cdot (x_j^{\otimes 4} - (x_j \otimes x_j) \tilde{\otimes} I + I \tilde{\otimes} I)(I, I, \alpha, \alpha) \\ &= (\sigma(w_i^\top x_j) + \xi_j) \cdot ((x^\top \alpha)^2 x^{\otimes 2} - (\alpha^\top x)^2 I - 2(\alpha^\top x)(x \alpha^\top + \alpha x^\top) - x x^\top + 2\alpha \alpha^\top + I), \end{aligned}$$

and $C_i(x_j) = \mathbb{1}(\|x_j\| \leq B) \cdot B_i(x_j)$. Then from Claim E.5 we have $\mathbb{E}[B_i(x_j)] = m_{4,i}(\bar{w}_i^\top \alpha)^2 \bar{w}_i \bar{w}_i^\top$. We calculate

$$\begin{aligned} &\sigma(w_i^\top x_j) \cdot (x_j^{\otimes 4} - (x_j \otimes x_j) \tilde{\otimes} I + I \tilde{\otimes} I)(I, I, \alpha, \alpha) \\ &\lesssim (|w_i^\top x_j|^{p+1} + |\phi(0)|) \cdot ((x_j^\top \alpha)^2 \|x_j\|^2 + 1 + \|x_j\|^2 + (\alpha^\top x_j)^2) \\ &\lesssim |w_i|^{p+1} \cdot |x_j|^{p+5}, \end{aligned}$$

By Assumption B.5, using Claim E.1 and $B \geq d \cdot \text{poly} \log(d)$ we have

$$\begin{aligned} \|\mathbb{E}[C_i(x_j)] - m_{4,i}(\bar{w}_i^\top \alpha)^2 \bar{w}_i \bar{w}_i^\top\| &\lesssim \mathbb{E}[\mathbb{1}_{\|x_j\| \geq B} |w_i|^{p+1} \cdot |x_j|^{p+5}] \\ &\lesssim (\|w_i\| d)^{p+5} \cdot e^{-\Omega(d \log d)} \\ &\lesssim \epsilon. \end{aligned}$$

Also, $\frac{1}{2}|m_{4,i}|(\bar{w}_i^\top \alpha)^2 \leq \|\mathbb{E}[C_i(x_j)]\| \leq 2|m_{4,i}|(\bar{w}_i^\top \alpha)^2$.

For any constant $t \in (0, 1)$, we have with probability $1 - t/4$,

$$\begin{aligned} \|C_i(x_j)\| &\lesssim (|w_i^\top x_j|^{p+1} + |\phi(0)| + |\xi_j|) \cdot ((x_j^\top \alpha)^2 \|x_j\|^2 + 1 + \|x_j\|^2 + (\alpha^\top x_j)^2) \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta) \cdot d \cdot \text{poly}(\log(d/t)) \end{aligned}$$

where the first step comes from Assumption B.5 and the second step comes from Claim E.2 and Claim E.3.

Using Claim E.4, we have

$$\begin{aligned} \|\mathbb{E}[C_i(x_j)^2]\| &\lesssim (\mathbb{E}[(\phi(w_i^\top x_j) + \xi_j)^4])^{1/2} (\mathbb{E}[(x_j^\top \alpha)^8])^{1/2} (\mathbb{E}[\|x_j\|^4])^{1/2} \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta)^2 d. \end{aligned}$$

Furthermore we have,

$$\max_{\|a\|=1} (\mathbb{E}[(a^\top C_i(x_j) a)^2])^{1/2} \lesssim (\mathbb{E}[(\phi(w_i^\top x_j) + \xi_j)^4])^{1/4} \lesssim \|w_i\|^{p+1} + |\phi(0)| + \vartheta.$$

Then by Claim E.6, with probability at least $1 - t$,

$$\begin{aligned} &\left\| m_{4,i}(\bar{w}_i^\top \alpha)^2 \bar{w}_i \bar{w}_i^\top - \frac{1}{|S|} \sum_{x_j \in S} C_i(x_j) \right\| \\ &\leq \left\| m_{4,i}(\bar{w}_i^\top \alpha)^2 \bar{w}_i \bar{w}_i^\top - \mathbb{E}[C_i(x_j)] \right\| + \left\| \mathbb{E}[C_i(x_j)] - \frac{1}{|S|} \sum_{x_j \in S} C_i(x_j) \right\| \\ &\lesssim \epsilon |m_{4,i}|(\bar{w}_i^\top \alpha)^2 + \epsilon. \end{aligned}$$

Summing up all components $i \in [k]$, we proved Eq (9). Eq (7) and Eq (8) can be shown similarly. \square

Lemma B.18. Let $V \in \mathbb{R}^{d \times k}$ be an orthogonal matrix. Let \widehat{R}_3 be computed in Line 20 of Algorithm 7 and $R_3 = P_3(V, V, V)$. Suppose

$$m_0 = \min_{i \in [k]} \{|m_{j_3, i}|^2 (\bar{w}_i^\top \alpha)^{2(j_3-3)}\}$$

and

$$|S| \gtrsim d \cdot \text{poly}(\kappa, \vartheta, \log(d/t)) / (\epsilon^2 m_0)$$

then with probability at least $1 - t$,

$$\|R_3 - \widehat{R}_3\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{j_3, i} (\bar{w}_i^\top \alpha)^{j_3-3}| + \epsilon.$$

Proof. From the definition of R_3 , it suffices to bound $\|M_3(V, V, V) - \widehat{M}_3(V, V, V)\|$ and $\|M_4(V, V, V, \alpha) - \widehat{M}_4(V, V, V, \alpha)\|$. The proof is similar to the previous one.

Specifically, we show that with probability at least $1 - t/4$,

$$\|M_3(V, V, V) - \widehat{M}_3(V, V, V)\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{3, i}| + \epsilon. \quad (10)$$

$$\|M_4(V, V, V, \alpha) - \widehat{M}_4(V, V, V, \alpha)\| \lesssim \epsilon \sum_{i=1}^k |v_i m_{4, i} (\bar{w}_i^\top \alpha)| + \epsilon. \quad (11)$$

Recall that for sample $(x_j, y_j) \in S$, $y_j = \sum_{i=1}^k v_i \sigma(w_i^\top x_j) + \xi_j$ where ξ_j is independent of x_j . Consider each component $i \in [k]$. Define $T_i(x_j), S_i(x_j) \in \mathbb{R}^{k \times k \times k}$:

$$\begin{aligned} T_i(x_j) &= (\sigma(w_i^\top x_j) + \xi_j) \\ &\quad \cdot (x_j^\top \alpha \cdot v(x) \otimes^3 - (V^\top \alpha) \otimes (v(x) \otimes v(x)) - \alpha^\top x \cdot v(x) \otimes I + (V^\top \alpha) \otimes I), \\ S_i(x_j) &= \mathbb{1}(\|x_j\| \leq B) \cdot T_i(x_j) \end{aligned}$$

where $v(x) = V^\top x$. Flatten $T_i(x_j)$ along the first dimension to obtain $B_i(x_j) \in \mathbb{R}^{k \times k^2}$, flatten $S_i(x_j)$ along the first dimension to obtain $C_i(x_j) \in \mathbb{R}^{k \times k^2}$.

From Claim E.7, $\mathbb{E}[B_i(x_j)] = m_{4, i} (\alpha^\top \bar{w}_i) (V^\top \bar{w}_i) \text{vec}((V^\top \bar{w}_i) (V^\top \bar{w}_i)^\top)^\top$. Therefore we have,

$$\|\mathbb{E}[B_i(x)]\| = |m_{4, i} (\alpha^\top \bar{w}_i)| \cdot \|V^\top \bar{w}_i\|^3.$$

We calculate

$$\begin{aligned} \|\mathbb{E}_{\xi_j}[B_i(x_j)]\| &\lesssim (|w_i^\top x_j|^{p+1} + |\phi(0)|) \cdot ((x_j^\top \alpha)^2 \|V^\top x_j\|^3 \\ &\quad + 3 \|V^\top x_j\|^3 + 3 |x_j^\top \alpha| \|V^\top x_j\| \sqrt{k} + 3 \|V^\top \alpha\| \sqrt{k}) \\ &\lesssim \sqrt{k} \cdot \|w_i\|^{p+1} \|x_j\|^{p+6}. \end{aligned}$$

By Assumption B.5, using Claim E.1 and $B \geq d \cdot \text{poly} \log(d)$,

$$\begin{aligned} &\|\mathbb{E}[C_i(x_j)] - m_{4, i} (\alpha^\top \bar{w}_i) (V^\top \bar{w}_i) \text{vec}((V^\top \bar{w}_i) (V^\top \bar{w}_i)^\top)^\top\| \\ &\lesssim \mathbb{E}[\mathbb{1}_{\|x_j\| \leq B} \sqrt{k} \|w_i\|^{p+1} \|x_j\|^{p+6}] \\ &\leq \epsilon. \end{aligned}$$

For any constant $t \in (0, 1)$, we have with probability $1 - t$,

$$\begin{aligned} \|C_i(x_j)\| &\lesssim (|w_i^\top x_j|^{p+1} + |\phi(0)| + |\xi_j|) \cdot ((x_j^\top \alpha)^2 \|V^\top x_j\|^3 \\ &\quad + 3 \|V^\top x_j\|^3 + 3 |x_j^\top \alpha| \|V^\top x_j\| \sqrt{k} + 3 \|V^\top \alpha\| \sqrt{k}) \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta) k^{3/2} \text{poly}(\log(d/t)) \end{aligned}$$

where the first step comes from Assumption B.5 and the second step comes from Claim E.2 and Claim E.3.

Using Claim E.4, we have

$$\begin{aligned} \|\mathbb{E}[C_i(x_j)C_i(x_j)^\top]\| &\lesssim (\mathbb{E}[(\phi(w_i^\top x_j) + \xi_j)^4])^{1/2} (\mathbb{E}[(\alpha^\top x_j)^4])^{1/2} (\mathbb{E}[\|V^\top x_j\|^6])^{1/2} \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta)^2 k^{3/2}. \end{aligned}$$

and

$$\begin{aligned} &\|\mathbb{E}[C_i(x_j)^\top C_i(x_j)]\| \\ &\lesssim (\mathbb{E}[(\phi(w_i^\top x_j) + \xi_j)^4])^{1/2} (\mathbb{E}[(\alpha^\top x_j)^4])^{1/2} (\mathbb{E}[\|V^\top x_j\|^4])^{1/2} \\ &\quad \cdot \left(\max_{\|A\|_F=1} \mathbb{E}[\langle A, (V^\top x_j)(V^\top x_j)^\top \rangle^4] \right)^{1/2} \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta)^2 k^2. \end{aligned}$$

Furthermore we have,

$$\begin{aligned} &\max_{\|a\|=\|b\|=1} (\mathbb{E}[(a^\top C_i(x_j)b)^2])^{1/2} \\ &\lesssim (\mathbb{E}[(\phi(w_i^\top x_j) + \xi_j)^4])^{1/4} (\mathbb{E}[(\alpha^\top x_j)^4])^{1/4} \max_{\|a\|=1} (\mathbb{E}[(a^\top V^\top x_j)^4])^{1/2} \\ &\quad \cdot \max_{\|A\|_F=1} (\mathbb{E}[\langle A, (V^\top x_j)(V^\top x_j)^\top \rangle^4])^{1/2} \\ &\lesssim (\|w_i\|^{p+1} + |\phi(0)| + \vartheta)k. \end{aligned}$$

Then by Claim E.6, with probability at least $1 - t$,

$$\begin{aligned} &\left\| m_{4,i}(\alpha^\top \bar{w}_i)(V^\top \bar{w}_i) \text{vec}((V^\top \bar{w}_i)(V^\top \bar{w}_i)^\top)^\top - \frac{1}{|S|} \sum_{x_j \in S} C_i(x_j) \right\| \\ &\leq \left\| m_{4,i}(\alpha^\top \bar{w}_i)(V^\top \bar{w}_i) \text{vec}((V^\top \bar{w}_i)(V^\top \bar{w}_i)^\top)^\top - \mathbb{E}[C_i(x_j)] \right\| \\ &\quad + \left\| \mathbb{E}[C_i(x_j)] - \frac{1}{|S|} \sum_{x_j \in S} C_i(x_j) \right\| \\ &\lesssim \epsilon |v_i m_{4,i}(\bar{w}_i^\top \alpha)| + \epsilon. \end{aligned}$$

Summing up all neurons $i \in [k]$, we proved Eq (11). Eq (10) can be shown similarly. \square

Lemma B.19. *Let \widehat{Q}_1 and \widehat{Q}_2 be computed in Line 28 of Algorithm 7. Let Q_1 be defined by Eq 4 and Q_2 be defined by Eq 5. Suppose*

$$m_0 = \min_{i \in [k]} \{|m_{j_1,i}|^2 (\bar{w}_i^\top \alpha)^{2(j_1-1)}, |m_{j_2,i}|^2 (\bar{w}_i^\top \alpha)^{2(j_2-2)}\}$$

and

$$|S| \gtrsim d \cdot \text{poly}(\kappa, \vartheta, \log(d/t)) / (\epsilon^2 m_0)$$

then with probability at least $1 - t$,

$$\begin{aligned} \|Q_1 - \widehat{Q}_1\| &\lesssim \epsilon \sum_{i=1}^k |v_i m_{j_1,i}(\bar{w}_i^\top \alpha)^{j_1-1}| + \epsilon, \\ \|Q_2 - \widehat{Q}_2\| &\lesssim \epsilon \sum_{i=1}^k |v_i m_{j_2,i}(\bar{w}_i^\top \alpha)^{j_2-2}| + \epsilon. \end{aligned}$$

Proof. Recall the expression of Q_1 and Q_2 ,

$$Q_1 = M_{I_1}(I, \underbrace{\alpha, \dots, \alpha}_{(j_1-1) \text{ } \alpha\text{'s}}) = \sum_{i=1}^k v_i c_{j_1} \|w_i\|^{p+1} (\alpha^\top \bar{w}_i)^{j_1-1} \bar{w}_i,$$

$$Q_2 = M_{j_2}(V, V, \underbrace{\alpha, \dots, \alpha}_{(j_2-2) \text{ } \alpha\text{'s}}) = \sum_{i=1}^k v_i c_{j_2} \|w_i\|^{p+1} (\alpha^\top \bar{w}_i)^{j_2-2} (V^\top \bar{w}_i) (V^\top \bar{w}_i)^\top.$$

The proof is essentially similar to Lemma B.17 and Lemma B.18. \square

We also use the following Lemmata from [47, 90].

Lemma B.20 (Adapted from Theorem 3 of [47]). *Given a tensor $\hat{T} = \sum_{i=1}^k \pi_i u_i^{\otimes 3} + \epsilon R \in \mathbb{R}^{d \times d \times d}$. Assume incoherence $u_i^\top u_j \leq \mu$. Let $L_0 := (\frac{50}{1-\mu^2})^2$ and $L \geq L_0 \log(15d(k-1)/t)^2$. Then there exists an algorithm such that, with probability at least $1-t$, for every u_i , the algorithm returns a \tilde{u}_i such that*

$$\|\tilde{u}_i - u_i\|_2 \leq O\left(\frac{\sqrt{\|\pi\|_1 \pi_{\max}}}{\pi_{\min}^2} \cdot \frac{\|V^\top\|_2^2}{1-\mu^2} \cdot (1+C(t))\right) \epsilon + o(\epsilon),$$

where $C(t) := \log(kd/t) \sqrt{d/L}$ and V is the inverse of the full-rank extension of $(u_1 \dots u_k)$ with unit-norm columns.

Lemma B.21 (Adapted from Lemma E.6 of [90]). *Let P_2 be defined as in Definition B.16 and \hat{P}_2 be its empirical version calculated in Line 4 of Algorithm 7. Let $U \in \mathbb{R}^{d \times k}$ be the orthogonal column span of $W \in \mathbb{R}^{d \times k}$. Assume $\|\hat{P}_2 - P_2\| \leq s_k(P_2)/10$. Let C be a large enough positive number such that $C > 2\|P_2\|$. Then after $T = O(\log(1/\epsilon))$ iterations, the $V \in \mathbb{R}^{d \times k}$ computed in Algorithm 7 will satisfy*

$$\|UU^\top - VV^\top\| \lesssim \|\hat{P}_2 - P_2\|/s_k(P_2) + \epsilon,$$

which implies

$$\|(I - VV^\top)w_i\| \lesssim (\|\hat{P}_2 - P_2\|/s_k(P_2) + \epsilon) \|w_i\|.$$

Lemma B.22 (Adapted from Lemma E.13 in [90]). *Let $U \in \mathbb{R}^{d \times k}$ be the orthogonal column span of W^* . Let $V \in \mathbb{R}^{d \times k}$ denote an orthogonal matrix satisfying that $\|VV^\top - UU^\top\| \leq \hat{\delta}_2 \lesssim 1/(\kappa^2 \sqrt{k})$. For each $i \in [k]$, let \hat{u}_i denote the vector satisfying $\|s_i \hat{u}_i - V^\top \bar{w}_i\| \leq \hat{\delta}_3 \lesssim 1/(\kappa^2 \sqrt{k})$. Let Q_1 be defined as in Eq (4) and \hat{Q}_1 be the empirical version of Q_1 such that $\|Q_1 - \hat{Q}_1\| \leq \hat{\delta}_4 \|Q_1\| \leq \frac{1}{4} \|Q_1\|$.*

Let $z^* \in \mathbb{R}^k$ and $\hat{z} \in \mathbb{R}^k$ be defined as in Eq (6) and Line 29. Then

$$\|\hat{z}_i - z_i^*\| \leq (\kappa^4 k^{3/2} (\hat{\delta}_2 + \hat{\delta}_3) + \kappa^2 k^{1/2} \hat{\delta}_4) \|z^*\|_1.$$

Lemma B.23 (Adapted from Lemma E.14 in [90]). *Let $U \in \mathbb{R}^{d \times k}$ be the orthogonal column span of W^* and V be an orthogonal matrix satisfying that $\|VV^\top - UU^\top\| \leq \hat{\delta}_2 \lesssim 1/(\kappa \sqrt{k})$. For each $i \in [k]$, let \hat{u}_i denote the vector satisfying $\|s_i \hat{u}_i - V^\top \bar{w}_i^*\| \leq \hat{\delta}_3 \lesssim 1/(\sqrt{k} \kappa^3)$.*

Let Q_2 be defined as in Eq (5) and \hat{Q}_2 be the empirical version of Q_2 such that $\|Q_2 - \hat{Q}_2\|_F \leq \hat{\delta}_4 \|Q_2\|_F \leq \frac{1}{4} \|Q_2\|_F$. Let $r^* \in \mathbb{R}^k$ and $\hat{r} \in \mathbb{R}^k$ be defined as in Eq (6) and Line 29. Then

$$\|\hat{r}_i - r_i^*\| \leq (k^3 \kappa^8 \hat{\delta}_3 + \kappa^2 k^2 \hat{\delta}_4) \|r^*\|.$$

Now we are in the position of proving Theorem B.13.

Proof. Consider Algorithm 7. First, by Lemma B.21 and Lemma B.17, we have

$$\begin{aligned} \|VV^\top \bar{w}_i - \bar{w}_i\| &\leq (\|\hat{P}_2 - P_2\|/s_k(P_2) + \epsilon) \\ &\leq (\text{poly}(k, \kappa) \|\hat{P}_2 - P_2\| + \epsilon) \\ &\leq \text{poly}(k, \kappa) \epsilon. \end{aligned} \tag{12}$$

Next, combining Lemma B.20 and Lemma B.18, we have

$$\|V^\top \bar{w}_i - s_i \hat{u}_i\| \leq \text{poly}(k, \kappa) \|\hat{R}_3 - R_3\| \leq \epsilon \text{poly}(k, \kappa). \quad (13)$$

It thus follows that

$$\begin{aligned} \|\bar{w}_i - s_i V \hat{u}_i\| &\leq \|VV^\top \bar{w}_i - \bar{w}_i\| + \|VV^\top \bar{w}_i - V s_i \hat{u}_i\| \\ &= \|VV^\top \bar{w}_i - \bar{w}_i\| + \|V^\top \bar{w}_i - s_i \hat{u}_i\| \\ &\leq \epsilon \text{poly}(k, \kappa), \end{aligned} \quad (14)$$

where the first step applies triangle inequality and the last step uses Eq (12) and Eq (13).

We proceed to bound the error in \hat{r} and \hat{z} . We have,

$$\begin{aligned} |\hat{r}_i - r_i^*| &\lesssim \text{poly}(k, \kappa) (\|Q_2 - \hat{Q}_2\| + \|s_i \hat{u}_i - V^\top \bar{w}_i\|) \cdot \|r^*\| \\ &\lesssim \epsilon \text{poly}(k, \kappa) \cdot \|r^*\| \end{aligned} \quad (15)$$

where the first step comes from Lemma B.23 and the second step comes from Lemma B.19 and Eq (14). Furthermore,

$$\begin{aligned} |\hat{z}_i - z_i^*| &\lesssim \text{poly}(k, \kappa) \cdot (\|Q_1 - \hat{Q}_1\| + \|s_i \hat{u}_i - V^\top \bar{w}_i\| + \|VV^\top - UU^\top\|) \cdot \|z^*\|_1 \\ &\lesssim \epsilon \text{poly}(k, \kappa) \|z^*\|_1, \end{aligned} \quad (16)$$

where the first step comes from Lemma B.22 and the second step comes from combining Lemma B.19, Lemma B.21, and Eq (14). Finally, combining Eq (15), Eq (16) and Eq (14), the output in Line 32 satisfies $\|\hat{w}_i - w_i\|_F \leq \epsilon \text{poly}(k, \kappa) \cdot \|w_i\|_F$. Since v_i are discrete values, they are exactly recovered. \square

B.5.2 Exact recovery of neural networks from noiseless data

In this section we prove Theorem B.15. Similar to Appendix B.5.1, denote $W = [w_1, \dots, w_k]^\top$ where $w_i \in \mathbb{R}^d$ and $\widehat{W} = [\hat{w}_1, \dots, \hat{w}_k]^\top$. We use \mathcal{D} to denote the distribution of $x \sim \mathcal{N}(0, I_d)$ and $y = \langle v, \sigma(Wx) \rangle$. We define the empirical loss for explored features and the population loss as follows,

$$L_n(\widehat{W}) = \frac{1}{2n} \sum_{(x,y) \in S^{(1)}} \left(\sum_{i=1}^k v_i \sigma(\hat{w}_i^\top x) - y \right)^2, \quad (17)$$

$$L(\widehat{W}) = \frac{1}{2} \mathbb{E}_{\mathcal{D}} \left[\left(\sum_{i=1}^k v_i \sigma(\hat{w}_i^\top x) - y \right)^2 \right]. \quad (18)$$

Algorithm 8 Using method of moments and gradient descent to recover neural network parameters

- 1: **procedure** NEURALNETRECOVERY($S = \{(x_i, y_i) : i \in [n]\}$)
 - 2: Let $S^{(1)} \leftarrow \{(x, y) \in S : \|x\| \leq B\}$
 - 3: Compute $(v, \widehat{W}) \leftarrow \text{NEURALNETNOISYRECOVERY}(S^{(1)})$
 - 4: Find $W^{(1)}$ as the global minimum of $L_n(\cdot)$, where $L_n(\cdot)$ is defined in Eq (17).
 - 5: **Return** $(v, W^{(1)})$
 - 6: **end procedure**
-

Definition B.24. Let s_i be the i -th singular value of W , $\lambda := \prod_{i=1}^k (s_i/s_k)$. Let $\tau = (3s_1/2)^{4p} / \min_{z \in [s_k/2, 3s_1/2]} \{\rho^2(z)\}$.

We use the follow results adapted from [90]. The only difference is that the rewards are potentially truncated if $\|x\| \geq B$, and due to $B = d \cdot \text{poly} \log(d)$ we can bound its difference between standard Gaussian in the same way as Appendix B.5.1.

Lemma B.25 (Concentration, adapted from Lemma D.11 in [90]). *Let samples size $n \geq \epsilon^{-2} d \tau \text{poly}(\log(d/t))$, then with probability at least $1 - t$,*

$$\|\nabla^2 L_n(W) - \nabla^2 L(W)\| \lesssim k s_1^{2p} \epsilon + \text{poly}(B_W, d) e^{-\Omega(d)}.$$

Lemma B.26 (Adapted from Lemma D.16 in [90]). *Assume activation $\sigma(\cdot)$ satisfies Assumption B.8 and Assumption B.7. Then for any $t \in (0, 1)$, if $n \geq d \cdot \text{poly}(\log(d/t))$, with probability at least $1 - t$, for any \widehat{W} (which is not necessarily to be independent of samples) satisfying $\|W - \widehat{W}\| \leq s_k/4$, we have*

$$\|\nabla^2 L_n(\widehat{W}) - \nabla^2 L_n(W)\| \leq k s_1^p \|W - \widehat{W}\| d^{(p+1)/2}.$$

Now we prove Theorem B.15.

Proof. The exact recovery consists of first finding (exact) v and (approximate) \widehat{W} close enough to W by tensor method (Appendix B.5.1), and then minimizing the empirical loss $L_n(\cdot)$. We will prove that $L_n(\cdot)$ is locally strongly convex, thus we find the precise W .

From Lemma D.3 from [90] we know:

$$\Omega(\rho(s_k)/\lambda)I \preceq \nabla^2 L(W) \preceq O(k s_1^{2p})I. \quad (19)$$

Combining Lemma B.25, $d \geq \log(B_W/\lambda)$, and $n \geq \frac{k^2 \lambda^2 s_1^{4p}}{\rho^2(s_k)} d \tau \text{poly}(\log(d/t))$, we know $\nabla^2 L_n(W)$ must be positive definite.

Next we uniformly bound Lipschitzness of $\nabla^2 L_n$. From Lemma B.26 there exists a universal constant c , such that for all \widehat{W} that satisfies $\|W - \widehat{W}\| \leq c k s_1^{2p} / (k s_1^p d^{(p+1)/2}) = c s_1^p d^{-(p+1)/2}$, $\nabla L_n^2(\widehat{W}) \gtrsim k s_1^{2p}$ holds uniformly. So there is a unique minimizer of L_n in this region.

Notice $L_n(W) = 0$, therefore we can find W by directly minimizing the empirical loss as long as we find any \widehat{W} in this region. This can be achieved by tensor method in Appendix B.5.1. We thus complete the proof. \square

C Omitted Proofs in Section 4

For the proofs of Theorem 4.2, Example 4.3, and Example 4.4, we refer the readers to [36].

Lemma C.1. *Consider the polynomial family $\mathcal{F}_\mathcal{V}$ of dimension D . Assume that $n > 2D$. For any $E \in \mathbb{R}^d$ that is of positive measure, by sampling n samples $\{x_i\}$ i.i.d. from $\mathbb{P}_{x \in \mathcal{N}(0, I_d)}(\cdot | x \in E)$ and observing the noiseless feedbacks $y_i = f^*(x_i)$, one can almost surely uniquely determine the f^* by solving the system of equations $y_i = f(x_i)$, $i = 1, \dots, n$, for $f \in \mathcal{F}_\mathcal{V}$.*

Proof. By Theorem 4.2, there exists a set $N \in \mathbb{R}^d \times \dots \times \mathbb{R}^d$ of Lebesgue measure zero, such that if $(x_1, \dots, x_n) \notin N$, one can uniquely determine the f^* by the observations on the n samples. Therefore, we only need to show that with probability 1, the sampling procedure returns $(x_1, \dots, x_n) \notin N$. This is because

$$\begin{aligned} \mathbb{P}(x_1, \dots, x_n \in N) &= \mathbb{P}_{x_i \in \mathcal{N}(0, I_d)}((x_1, \dots, x_n) \in N \mid x_1, \dots, x_n \in E) \\ &= \frac{\mathbb{P}_{x_i \in \mathcal{N}(0, I_d)}((x_1, \dots, x_n) \in N \cap (E \times \dots \times E))}{\mathbb{P}_{x_i \in \mathcal{N}(0, I_d)}((x_1, \dots, x_n) \in (E \times \dots \times E))} \\ &= \frac{0}{[\mathbb{P}_{x_1 \in \mathcal{N}(0, I_d)}(x_1 \in E)]^n} \\ &= 0. \end{aligned}$$

\square

By Lemma C.1 above, it is not hard to see that Algorithms 4 and 5 work.

D Omitted Constructions and Proofs in Subsection 4.1

Construction of the Reward Functions The following construction of the polynomial hard case is adopted from [36].

Let d be the dimension of the feature space. Let e_i denotes the i -th standard orthonormal basis of \mathbb{R}^d , i.e., e_i has only one 1 at the i -th entry and 0's for other entries. Let p denote the highest order of the polynomial. We assume $d \gg p$. We use Λ to denote a subset of the p -th multi-indices

$$\Lambda = \{(\alpha_1, \dots, \alpha_p) | 1 \leq \alpha_1 \leq \dots \leq \alpha_p \leq d\}.$$

For an $\alpha = (\alpha_1, \dots, \alpha_p) \in \Lambda$, denote $M_\alpha = e_{\alpha_1} \otimes \dots \otimes e_{\alpha_p}$, $x_\alpha = e_{\alpha_1} + \dots + e_{\alpha_p}$.

The model space \mathcal{M} is a subset of rank-1 p -th order tensors, which is defined as $\mathcal{M} = \{M_\alpha | \alpha \in \Lambda\}$. We define two subsets of feature space \mathcal{F}_0 and \mathcal{F} as $\mathcal{F}_0 = \{x_\alpha | \alpha \in \Lambda\}$, $\mathcal{F} = \text{conv}(\mathcal{F}_0)$. For $M_\alpha \in \mathcal{M}$, $x \in \mathcal{F}$, define $r(M_\alpha, x)$ as $r(M_\alpha, x) = \langle M_\alpha, x^{\otimes p} \rangle = \prod_{i=1}^p \langle e_{\alpha_i}, x \rangle$. We assume that for each level h , there is a $M^{(h)} = M_{\alpha^{(h)}} \in \mathcal{M}$, and the noiseless reward is $r_h(s, a) = r(M^{(h)}, \phi_h(s, a))$.

We have the following properties.

Proposition D.1 ([36]). *For $M_\alpha \in \mathcal{M}$ and $x_{\alpha'} \in \mathcal{F}_0$, we have*

$$r(M_\alpha, x_{\alpha'}) = \mathbb{I}_{\{\alpha=\alpha'\}}.$$

Proposition D.2. *For $M_\alpha \in \mathcal{M}$, we have*

$$\max_{x \in \mathcal{F}} r(M_\alpha, x) = 1.$$

proof of Proposition D.2. For all $x \in \mathcal{F}$, since $\mathcal{F} = \text{conv}(\mathcal{F}_0)$, we can write

$$x = \sum_{\alpha \in \Lambda} p_\alpha (e_{\alpha_1} + \dots + e_{\alpha_p}),$$

where $\sum_{\alpha \in \Lambda} p_\alpha = 1$ and $p_\alpha \geq 0$. Therefore,

$$r(M_{\alpha'}, x) = \prod_{i=1}^p \langle e_{\alpha'_i}, x \rangle.$$

Plug in the expression of x , we have

$$\begin{aligned} \langle e_{\alpha'_i}, x \rangle &= \sum_{\alpha} p_\alpha \langle e_{\alpha'_i}, e_{\alpha_1} + \dots + e_{\alpha_p} \rangle \\ &= \sum_{\alpha} p_\alpha \mathbb{I}_{\{\alpha'_i \in \alpha\}} \\ &\leq \sum_{\alpha} p_\alpha = 1. \end{aligned}$$

Therefore,

$$\begin{aligned} r(M_{\alpha'}, x) &= \prod_{i=1}^p \langle e_{\alpha'_i}, x \rangle \\ &= \left(\sum_{\alpha} p_\alpha \mathbb{I}_{\{\alpha'_1 \in \alpha\}} \right) \cdots \left(\sum_{\alpha} p_\alpha \mathbb{I}_{\{\alpha'_p \in \alpha\}} \right) \\ &\leq 1. \end{aligned}$$

Finally, since $r(M_{\alpha'}, x_{\alpha'}) = 1$, we have $\max_{x \in \mathcal{F}} r(M_\alpha, x) = 1$. \square

MDP constructions Consider a family of MDPs with only two states $\mathcal{S} = \{S_{\text{good}}, S_{\text{bad}}\}$. The action set \mathcal{A} is set to be \mathcal{F} . Let f be a mapping from \mathcal{F} to \mathcal{F}_0 such that f is identity when restricted to \mathcal{F}_0 . For all level $h \in [H]$, we define the feature map $\phi_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{F}$ to be

$$\phi_h(s, a) = \begin{cases} a & \text{if } s = S_{\text{good}}, \\ f(a) & \text{if } s = S_{\text{bad}}. \end{cases}$$

Given an unknown sequence of indices $\alpha^{(1)}, \dots, \alpha^{(H)}$, the reward function at level h is $r_h(s, a) = r(M_{\alpha^{(h)}}, \phi_h(s, a))$. Specifically, we have

$$r_h(S_{\text{good}}, a) = r(M_{\alpha^{(h)}}, a), \quad r_h(S_{\text{bad}}, a) = r(M_{\alpha^{(h)}}, f(a)).$$

The transition P_h is constructed as

$$P_h(S_{\text{bad}}|s, a) = 1 \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}.$$

This construction means it is impossible for the online scenarios to reach the good state for $h > 1$.

The next proposition shows that Q_h^* is polynomial realizable and falls into the case of Example 4.4.

Proposition D.3. *We have for all $h \in [H]$ and $s \in \mathcal{S}, a \in \mathcal{A}$, $V_h^*(s) = H - h + 1$ and $Q_h^*(s, a) = r_h(s, a) + H - h + 1$. Furthermore, $Q_h^*(s, a)$, viewed as the function of $\phi_h(s, a)$, is a polynomial of the form $q_h(U_h \phi_h(s, a))$ for some degree- p polynomial q_h and $U_h \in \mathbb{R}^{p \times d}$.*

proof of Proposition D.3. First notice that by Proposition D.2, for all $h \in [H]$ and $s \in \mathcal{S}$, we have

$$\max_{a \in \mathcal{A}} r_h(s, a) = 1.$$

Therefore, by induction, suppose we have proved for all $s', V_{h+1}^*(s') = H - h$, then we have

$$\begin{aligned} V_h^*(s) &= \max_{a \in \mathcal{A}} Q_h^*(s, a) \\ &= \max_{a \in \mathcal{A}} \{r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot|s, a)} [V_{h+1}^*(s')]\} \\ &= 1 + H - h. \end{aligned}$$

Then we have $Q_h^*(s, a) = r_h(s, a) + H - h + 1$.

Furthermore, we have

$$\begin{aligned} Q_h^*(s, a) &= r_h(s, a) + H - h + 1 \\ &= r(M_{\alpha^{(h)}}, \phi_h(s, a)) + H - h + 1 \\ &= \prod_{i=1}^p \langle e_{\alpha_i^{(h)}}, \phi_h(s, a) \rangle + H - h + 1 \\ &= q_h(U_h \phi_h(s, a)), \end{aligned}$$

where $q_h(x_1, \dots, x_p) = x_1 x_2 \cdots x_p + (H - h + 1)$ and $U_h \in \mathbb{R}^{p \times d}$ is a matrix with $e_{\alpha_i^{(h)}}$ as the i -th row. \square

Theorem D.4. *Under the online RL setting, any algorithm needs to play at least $\binom{d}{p} - 1 = \Omega(d^p)$ episodes to identify $\alpha^{(2)}, \dots, \alpha^{(H)}$ and thus to identify the optimal policy.*

proof of Theorem D.4. Under the online RL setting, any algorithm enters and remains in S_{bad} for $h > 1$. When $s_h = S_{\text{bad}}$, no matter what a_h the algorithm chooses, we have $\phi_h(s_h, a_h) = f(a_h) \in \mathcal{F}_0$. Notice that for any $M_{\alpha^{(h)}} \in \mathcal{M}$ and any $x_\alpha \in \mathcal{F}_0$, we have $r(M_{\alpha^{(h)}}, x_\alpha) = \mathbb{I}_{\{\alpha = \alpha^{(h)}\}}$ as Proposition D.1 suggests. Hence, we need to play $\binom{d}{p} - 1$ times at level h in the worst case to find out $\alpha^{(h)}$. The argument holds for all $h = 2, 3, \dots, H$. \square

Theorem D.5. *Under the generative model setting, by querying $2d(p+1)^p H = O(dH)$ samples, we can almost surely identify $\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(H)}$ and thus identify the optimal policy.*

proof of Theorem D.5. By Proposition D.3, we know that $Q_h^*(s, a)$, viewed as the function of $\phi_h(s, a)$, falls into the case of Example 4.4 with $k = p$.

Next, notice that for all $h \in [H]$, $\{\phi_h(s, a) \mid s \in \mathcal{S}, a \in \mathcal{A}\} = \mathcal{F}$. Although \mathcal{F} is not of positive measure, we can actually know the value of Q_h^* when $\phi_h(s, a)$ is in $\text{conv}(\mathcal{F}, \mathbf{0})$ since the reward is p -homogenous. Specifically, for every feature of the form $c \cdot \phi_h(s, a)$, where $0 \leq c \leq 1$ and $\phi_h(s, a) \in \mathcal{F}$, the reward is c^p times the reward of (s, a) . Therefore, to get the reward at $c \cdot \phi_h(s, a)$, we only need to query the generative model at (s, a) of level h , and then multiply the reward by c^p .

Notice that $\text{conv}(\mathcal{F}, \mathbf{0})$ is of positive Lebesgue measure. By Theorem 4.7, we know that only $2d(p+1)^p H = O(dH)$ samples are needed to determine the optimal policy almost surely. \square

E Technical claims

Claim E.1. Let $\chi^2(d)$ denote χ^2 -distribution with degree of freedom d . For any $t > 0$ we have,

$$\Pr_{z \sim \chi^2(d)} (z \geq d + 2t + 2\sqrt{dt}) \leq e^{-t}$$

We use the following facts from [90].

Claim E.2. Given a fixed vector $z \in \mathbb{R}^d$, for any $C \geq 1$ and $n \geq 1$, we have

$$\Pr_{x \sim \mathcal{N}(0, I_d)} [|\langle x, z \rangle|^2 \leq 5C\|z\|^2 \log n] \geq 1 - 1/(nd^C).$$

Claim E.3. For any $C \geq 1$ and $n \geq 1$, we have

$$\Pr_{x \sim \mathcal{N}(0, I_d)} [\|x\|^2 \leq 5Cd \log n] \geq 1 - 1/(nd^C).$$

Claim E.4. Let $a, b, c \geq 0$ be three constants, let $u, v, w \in \mathbb{R}^d$ be three vectors, we have

$$\mathbb{E}_{x \sim \mathcal{N}(0, I_d)} [u^\top x|^a |v^\top x|^b |w^\top x|^c] \approx \|u\|^a \|v\|^b \|w\|^c.$$

Claim E.5. Let $M_j, j \in [4]$ be defined in Definition B.11. For each $j \in [4]$, $M_j = \sum_{i=1}^k v_i m_{j,i} \bar{w}_i^{\otimes j}$.

Claim E.6. Let \mathcal{B} denote a distribution over $\mathbb{R}^{d_1 \times d_2}$. Let $d = d_1 + d_2$. Let B_1, B_2, \dots, B_n be i.i.d. random matrices sampled from \mathcal{B} . Let $\bar{B} = \mathbb{E}_{B \sim \mathcal{B}}[B]$ and $\hat{B} = \frac{1}{n} \sum_{i=1}^n B_i$. For parameters $m \geq 0, \gamma \in (0, 1), \nu > 0, L > 0$, if the distribution \mathcal{B} satisfies the following four properties,

- (1) $\Pr_{B \sim \mathcal{B}} [\|B\| \leq m] \geq 1 - \gamma;$
- (2) $\left\| \mathbb{E}_{B \sim \mathcal{B}}[B] \right\| > 0;$
- (3) $\max \left(\left\| \mathbb{E}_{B \sim \mathcal{B}}[BB^\top] \right\|, \left\| \mathbb{E}_{B \sim \mathcal{B}}[B^\top B] \right\| \right) \leq \nu;$
- (4) $\max_{\|a\|=\|b\|=1} \left(\mathbb{E}_{B \sim \mathcal{B}} \left[(a^\top B b)^2 \right] \right)^{1/2} \leq L.$

Then we have for any $0 < \epsilon < 1$ and $t \geq 1$, if

$$n \geq (18t \log d) \cdot (\nu + \|\bar{B}\|^2 + m\|\bar{B}\|\epsilon) / (\epsilon^2 \|\bar{B}\|^2) \quad \text{and} \quad \gamma \leq (\epsilon \|\bar{B}\| / (2L))^2$$

with probability at least $1 - 1/d^{2t} - n\gamma$,

$$\|\hat{B} - \bar{B}\| \leq \epsilon \|\bar{B}\|.$$

Claim E.7. Let P_2 and P_3 be defined in Definition B.16. Then

$$P_2 = \sum_{i=1}^k v_i m_{j_2, i} (\alpha^\top \bar{w}_i)^{j_2 - 2} \bar{w}_i^{\otimes 2}$$

and

$$P_3 = \sum_{i=1}^k v_i m_{j_3, i} (\alpha^\top \bar{w}_i)^{j_3 - 3} \bar{w}_i^{\otimes 3}.$$