# Provable General Function Class Representation Learning in Multitask Bandits and MDPs

## Rui Lu<sup>1</sup>, Andrew Zhao<sup>1</sup>, Simon S. Du<sup>2</sup>, Gao Huang<sup>1</sup>

<sup>1</sup>Department of Automation, BNRist, Tsinghua University
<sup>2</sup>Paul G. Allen School of Computer Science and Engineering, University of Washington {r-lu21,zqc21}@mails.tsinghua.edu.cn
ssdu@cs.washington.com, gaohuang@tsinghua.edu.cn

#### **Abstract**

While multitask representation learning has become a popular approach in reinforcement learning (RL) to boost the sample efficiency, the theoretical understanding of why and how it works is still limited. Most previous analytical works could only assume that the representation function is already known to the agent or from linear function class, since analyzing general function class representation encounters non-trivial technical obstacles such as generalization guarantee, formulation of confidence bound in abstract function space, etc. However, linear-case analysis heavily relies on the particularity of linear function class, while real-world practice usually adopts general non-linear representation functions like neural networks. This significantly reduces its applicability. In this work, we extend the analysis to general function class representations. Specifically, we consider an agent playing M contextual bandits (or MDPs) concurrently and extracting a shared representation function  $\phi$  from a specific function class  $\Phi$  using our proposed Generalized Functional Upper Confidence Bound algorithm (GFUCB). We theoretically validate the benefit of multitask representation learning within general function class for bandits and linear MDP for the first time. Lastly, we conduct experiments to demonstrate the effectiveness of our algorithm with neural net representation.

# 1 Introduction

Recently, reinforcement learning (RL) has achieved many successful applications in games [6, 34], robotics [23], and many other fields. However, due to the large cardinality of state space or action space in real-world problems, the large sample complexity has been a major problem for employing these RL algorithms in reality. A popular method called multitask representation learning tries to tackle this problem by extracting a shared low-dimensional representation function among multiple related tasks, then using a simple function (e.g., linear) on top of this common representation to solve each task[4, 7, 24].

Despite the empirical success for multitask representation learning, particularly in reinforcement learning because of its effectiveness in reducing sample complexity, the theoretical understanding about it is still limited. A march of works[37, 36, 22, 33, 25, 31, 16, 3, 11, 9, 41, 30] give results on function approximation in bandits and RL, which permits a representation. In these frameworks, an agent is considered playing M related tasks concurrently. Each task is a distinct contextual bandit or linear MDP problem  $^1$ , and all these M tasks share a common representation  $\phi \in \Phi$  where  $\Phi = \{\phi: \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^k\}$  is representation function class extracting a k-dimensional representation

<sup>&</sup>lt;sup>1</sup>Although the name of linear MDP contains term "linear", it actually has infinite degrees of freedom because the representation function  $\phi$  could be general non-linear function.

vector from state-action pair. Such representation function can reduce the complexity of problem from a huge space  $\mathcal{S} \times \mathcal{A}$  to a simple regression problem in k-dimensional space. The value approximation function class is defined by  $\mathcal{F} = \mathcal{L} \circ \Phi$ , here  $\circ$  means composition and  $\mathcal{L}$  means linear function, which means the value of any state-action pair (s,a) is linear in its representation  $\phi(s,a)$ .

However, previous analyses either assume  $\Phi$  is linear [40], or assume that the agent already knows the concrete function  $\phi$  [17, 21], which equivalently reduces to learning linear weight parameters. This limits their applicability, since general non-linear value estimation is ubiquitous and is the essence for the success of multitask representation learning. For instance, DQN[29] achieves great success by employing a deep network to approximate Q-value function. Also, assuming the agent already knows a good representation function is unrealistic in practice. Therefore, we aim to extend the analysis to unknown general non-linear representation functions. This would not only reveal the more essential benefit of multitask representation learning, but also inspire and facilitate future practice.

#### 1.1 Our Contribution

The focus of previous works on linear analyses has its own reasons. The particularity of linear function could circumvent many non-trivial obstacles in analysis, which hinders previous work from from being applicable to real world scenarios. For instance, the formulation of confidence set in linear parameter space is simply an ellipsoid, and its update is straightforward via covariance matrix. More importantly, linear function class generically ensures generalization. The analysis [18, 43, 26, 21] only requires the samples to span the whole input space to let the covariance matrix converge, then is able to derive uniform prediction error guarantee for the whole input space. However, generalization issue is much more complicated for general non-linear scenarios.

In summary, our work embraces following contributions, which solves the challenges for previous works and extends the analysis for the role of representation function in more general setting.

Eliminate the Dependency on Linearity. Towards general function class analysis, we adopt the idea of confidence set [32, 17]. The algorithm extends the idea of upper confidence bound and maintains a *confidence set* for all the possible value estimation functions. The confidence set contains all the functions whose total empirical error at step t is less than a predetermined bound  $\beta_t$ . As more seen data reveals more information about the environment, the confidence set will gradually shrink until converge. Therefore, our algorithm and analysis framework is applicable to general function class.

Note that designing  $\beta_t$  to achieve low regret for general function class  $\Phi$  is non-trivial. We firstly determine the concrete UCB form for general function class  $\beta_t(\Phi)$  and propose a straightforward algorithm called Generalized Functional Upper Confidence Bound (or GFUCB in abbreviation) for general non-linear function class approximation. We use Eluder dimension[32] to measure the complexity of the function class  $\Phi$  to give an efficient sample complexity that ensures generalization.

**Multihead Function Class.** To derive sharp regret bound for our algorithm and theoretically demonstrate the benefit of multitask representation learning, we firstly introduce multihead function class  $\mathcal{F}^{\otimes M}$ , which is the key technical contribution of our work. The efficacy of multitask representation learning essentially originates from the shared knowledge and structure among tasks. Hence it is vital and necessary to characterize such relation between multiples tasks that the agent simultaneously learns. However, such structure is absent in previous single task work [39, 32], and it calls for special techniques to analyze the efficiency for learning these correlated functions.

To this end, we introduce multihead function class, namely  $\mathcal{F}^{\otimes M}$  in section 4. This abstract function space captures the relation between different task functions, which concatenate the values of (s,a) for all M tasks together as the output. Being more compact by sharing a common backbone  $\phi$ , function in  $\mathcal{F}^{\otimes M}$  requires much fewer samples to learn compared to M independent tasks space  $\mathcal{F}^M$ . All the tasks contribute to shape a good representation, then feedback to each task for faster convergence. We formally prove that our algorithm enjoys regret bound as  $\tilde{O}\left(\sqrt{MT\dim_E(\mathcal{F})(Mk+\log\mathcal{N}(\Phi))}\right)$ , where T is the number of steps, M is the number of tasks and  $\mathcal{N}(\Phi)$  means the covering number of function space  $\Phi$ . We also extend the algorithm and analysis to multitask episodic RL with general value approximation under low inherent Bellman error. By simultaneously solving M different but correlated MDP tasks, our method is sample-efficient with

regret  $\tilde{O}\left(\sqrt{MTH\dim_E(\mathcal{F})(Mk+\log\mathcal{N}(\Phi)+MTH\mathcal{I}^2)}\right)$  where T is the number of episodes, H is planning horizon and  $\mathcal{I}$  denotes the inherent Bellman error.

To the best of our knowledge, this is the first provably sample efficient algorithm for general representation function bandits and linear MDP. It is comparable to the most optimal regret bound when  $\Phi$  is specialized to linear representation, and is better than the bounds which solve each task independently. This also theoretically explains how multitask representation learning reduces sample complexity. Essentially, the joint training for the shared representation function helps accelerate the convergence of the common backbone by having more samples from all the tasks.

**Empirical Value.** Finally, we conduct experiments to verify our theoretical result. We design a neural network based bandit environment and implement the GFUCB algorithm. Experimental results corroborate the effect of multitask representation learning in boosting sample efficiency in non-linear bandits. For the first time, the efficacy of the general representation algorithm proposed in theoretical analysis is validated in a proof-of-concept experiment.

#### 2 Related Work

In the supervised learning setting, a line of works have been done on multitask learning and representation learning with various assumptions [4, 15, 2, 5, 27, 8, 28, 14, 38]. These results assumed that all tasks share a joint representation function. It is also worth mentioning that [38] gave the method-of-moments estimator and built the confidence ball for the feature extractor, which inspired our algorithm for the infinite-action setting.

The benefit of representation learning has been studied in sequential decision-making problems, especially in RL domains. Arora et al. [3] proved that representation learning could reduce the sample complexity of imitation learning. D'eramo et al. [11] showed that representation learning could improve the convergence rate of the value iteration algorithm. Both require a probabilistic assumption similar to that in [28], and the statistical rates are of similar forms as those in [28]. Following these works, we study a special class of MDP called Linear MDP. Linear MDP [42, 21] is a popular model in RL, which uses linear function approximation to generalize large state-action space. [44] extends the definition to low inherent Bellman error (or IBE in short) MDPs. This model assumes that both the transition and the reward are near-linear in given features.

Recently, Yang et al. [40] showed multitask representation learning reduces the regret in linear bandits, using the framework developed by Du et al. [14]. Moreover, some works [17, 26, 21] proved results on the benefit of multitask representation learning RL with generative model or linear representation function. However, these works either restrict the representation function class to be linear, or the representation function is known to agent. This is unrealistic in real world practice, which limits these works' meaning.

The most relevant works that need to be mentioned is general function class value approximation for bandits and MDPs. Russo et al. [32] first proposed the concept of eluder dimension to measure the complexity of a function class and gave a regret bound for general function bandits using this dimension. Wang et al. [39] further proved that it can also be adopted in MDP problems. Dong et al. [12] extended the analysis with sequential Rademacher complexity. Inspired by these works, we adopt eluder dimension and develop our own analysis. But it should be pointed out that all those works focus on single task setting, which give a provable bound for just one single MDP or bandit problem. They lack the insight for why simultaneously dealing with multiple distinct but correlated tasks is more sample efficient. Our work aim to establish a framework to explain this. By considering locating the ground truth value function in multihead function space  $\mathcal{F}^{\otimes M}$  (see detailed definition in section 4), we are able to theoretically explain the main reason for the boost of sample efficiency. Informally speaking, the shared feature extraction backbone  $\phi$  receives samples from all the tasks, therefore accelerating the convergence for every single task compare with solving them separately.

#### 3 Preliminaries

#### 3.1 Notations

We use [n] to denote the set  $\{1,2,\ldots,n\}$  and  $\langle\cdot,\cdot\rangle$  to denote the inner product between two vectors. We use f(x)=O(g(x)) to represent  $f(x)\leq C\cdot g(x)$  holds for any  $x>x_0$  with some C>0 and  $x_0>0$ . Ignoring the logarithm term, we use  $f(x)=\tilde{O}(g(x))$ .

#### 3.2 Multitask Contextual Bandits

We first study multitask representation learning in contextual bandits. Each task  $i \in [M]$  is associated with an unknown function  $f^{(i)} \in \mathcal{F}$  from certain function class  $\mathcal{F}$ . At each step  $t \in [T]$ , the agent is given a context vector  $C_{t,i}$  from certain context space  $\mathcal{C}$  and a set of actions  $\mathcal{A}_{t,i}$  selected from certain action space  $\mathcal{A}$  for each task i. The agent needs to choose one action  $A_{t,i} \in \mathcal{A}_{t,i}$ , and then receives a reward as  $R_{t,i} = f^{(i)}(C_{t,i}, A_{t,i}) + \eta_{t,i}$ , where  $\eta_{t,i}$  is the random noise sampled from some i.i.d. distribution. The agent's goal is to understand function  $f^{(i)}$  and maximize the cumulative reward, or equivalently, minimize the total regret from all M tasks in T steps defined as below.

$$\operatorname{Reg}(T) \stackrel{\text{def}}{=} \sum_{t=1}^{T} \sum_{i=1}^{M} \left( f^{(i)}(C_{t,i}, A_{t,i}^{\star}) - f^{(i)}(C_{t,i}, A_{t,i}) \right),$$

where  $A_{t,i}^{\star} = \arg \max_{A \in \mathcal{A}_{t,i}} f^{(i)}(C_{t,i}, A)$  is the optimal action with respect to context  $C_{t,i}$  in task i.

#### 3.3 Multitask MDP

Going beyond contextual bandits, we also study how this shared low-dimensional representation could benefit the sequential decision making problem like Markov Decision Process (MDP). In this work, we study undiscounted episodic finite horizon MDP problem. Consider an MDP  $\mathcal{M}=(\mathcal{S},\mathcal{A},\mathcal{P},r,H)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\mathcal{P}$  is the transition dynamics,  $r(\cdot,\cdot)$  is the reward function and H is the planning horizon. The agent starts from an initial state  $s_1$  which can be either fixed or sampled from a certain distribution, then interacts with environment for H rounds. In the single task framework, at each round (also called level) h, the agent needs to perform an action  $a_h$  according to a policy function  $a_h = \pi_h(s_h)$ . Then the agent will receive a reward  $R_h(s_h,a_h)=r(s_h,a_h)+\eta_h$  where  $\eta_h$  again is the noise term. The environment then transits the state from  $s_h$  to  $s_{h+1}$  according to distribution  $\mathcal{P}(\cdot|s_h,a_h)$ . The estimation for action value function given following action policy  $\pi$  is defined as  $Q_h^\pi(s_h,a_h)=r(s_h,a_h)+\mathbb{E}\left[\sum_{t=h+1}^H R_t(s_t,\pi_t(s_t))\right]$ , and state value function is defined as  $V_h^\pi(s_h)=Q_h^\pi(s_h,\pi_h(s_h))$ . Note that there always exists a deterministic optimal policy  $\pi^\star$  for which  $V_h^{\pi^\star}(s)=\max_{x}V_h^\pi(s)$  and  $Q_h^{\pi^\star}(s,a)=\max_{x}Q_h^\pi(s,a)$ , we will denote them as  $V_h^\star(s)$  and  $Q_h^\star(s,a)$  for simplicity.

In the multitask setting, the agent gets a batch of states  $\{s_{h,t}^{(i)}\}_{i=1}^{M}$  simultaneously from M different MDP tasks  $\{\mathcal{M}^{(i)}\}_{i=1}^{M}$  at each round h in episode t, then performs a batch of actions  $\{\pi_t^i(s_{h,t}^{(i)})\}_{i=1}^{M}$  for each task  $i \in [M]$ . Every H rounds form an episode, and the agent will interact with the environment for totally T episodes. The goal for the agent is minimizing the regret defined as

$$\operatorname{Reg}(T) = \sum_{t=1}^{T} \sum_{i=1}^{M} V_1^{(i)\star} \left( s_{1,t}^{(i)} \right) - V_1^{\pi_t^i} \left( s_{1,t}^{(i)} \right),$$

where  $V_1^{(i)\star}$  is the optimal value of task i and  $s_{1,t}^{(i)}$  is the initial state for task i at episode t.

To let representation function play a role, it is assumed that all tasks share the same state space  $\mathcal S$  and action space  $\mathcal A$ . Moreover, there exists a representation function  $\phi: \mathcal S \times \mathcal A \mapsto \mathbb R^k$  such that action and state value function of all tasks  $\mathcal M^{(i)}$  is always (approximately) linear in this representation. For example, given a representation function  $\phi$ , the action value approximation function at level h is parametrized by a vector  $\boldsymbol{\theta}_h \in \mathbb R^k$  as  $Q_h[\phi, \boldsymbol{\theta}_h] \stackrel{\text{def}}{=} \langle \phi(s,a), \boldsymbol{\theta}_h \rangle$ , similar for  $V_h[\phi, \boldsymbol{\theta}_h](s) \stackrel{\text{def}}{=} \max_a \langle \phi(s,a), \boldsymbol{\theta}_h \rangle$ . We denote all such action value functions as  $\mathcal Q_h = \{Q_h[\phi, \boldsymbol{\theta}_h]: \phi \in \Phi, \boldsymbol{\theta}_h \in \mathbb R^k\}$ , also value function approximation space as  $\mathcal V_h = \{V_h[\phi, \boldsymbol{\theta}_h]: \phi \in \Phi, \boldsymbol{\theta}_h \in \mathbb R^k\}$ . Each task

 $\mathcal{M}^{(i)}$  is a linear MDP, which means  $\mathcal{Q}_h$  is always approximately close under Bellman operator  $\mathcal{T}_h(Q_{h+1})(s,a) \stackrel{\text{def}}{=} r_h(s,a) + \mathbb{E}_{s' \sim \mathcal{P}_h(\cdot|s,a)} \max_{a'} Q_{h+1}(s',a').$ 

**Linear MDP Definition.** A finite horizon MDP  $\mathcal{M} = (S, A, \mathcal{P}, r, H)$  is a linear MDP, if there exists a representation function  $\phi : S \times A \mapsto \mathbb{R}^k$  and its induced value approximation function class  $Q_h, h \in [H]$ , such that the inherent Bellman error[44]

$$\mathcal{I}_{h} \stackrel{\text{def}}{=} \sup_{Q_{h+1}} \inf_{Q_{h} \in \mathcal{Q}_{h}} \sup_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \left( Q_{h} - \mathcal{T}_{h} \left( Q_{h+1} \right) \right) \left( s, a \right) \right|,$$

is always smaller than some small constant  $\mathcal{I}$ .

The definition essentially assumes that for any Q-value approximation function  $Q_{h+1} \in \mathcal{Q}_{h+1}$  at level h+1, the Q-value function  $Q_h$  at level h induced by it can always be closely approximated in class  $\mathcal{Q}_h$ , which assures the accuracy through sequential levels.

#### 3.4 Eluder Dimension

To measure the complexity of a general function class f, we adopt the concept of eluder dimension [32]. First, define  $\epsilon$ -dependence and independence.

**Definition 1** ( $\epsilon$ -dependent). An input x is  $\epsilon$ -dependent on set  $X = \{x_1, x_2, \dots, x_n\}$  with respect to function class  $\mathcal{F}$ , if any pair of functions  $f, \tilde{f} \in \mathcal{F}$  satisfying  $\sqrt{\sum_{i=1}^n (f(x_i) - \tilde{f}(x_i))^2} \le \epsilon$  also satisfies  $|f(x) - \tilde{f}(x)| \le \epsilon$ . Otherwise, we call action x to be  $\epsilon$ -independent of data set X.

Intuitively,  $\epsilon$ -dependence captures the exhaustion of interpolation flexibility for function class  $\mathcal{F}$ . Given an unknown function f's value on set  $X = \{x_1, x_2, \dots, x_n\}$ , we are able to pin down its value on some particular input x with only  $\epsilon$ -scale prediction error.

**Definition 2** ( $\epsilon$ -eluder dimension). The  $\epsilon$ -eluder dimension  $\dim_E(\mathcal{F}, \epsilon)$  is the maximum length for a sequence of inputs  $x_1, x_2, \dots x_d \in \mathcal{X}$ , such that for some  $\epsilon' \geq \epsilon$ , every element is  $\epsilon'$ -independent of its predecessors.

This definition is similar to the definition of the dimensionality of a linear space, which is the maximum length of a sequence of vectors such that each one is linear independent to its predecessors. For instance, if  $\mathcal{F} = \{f(x) : \mathbb{R}^d \mapsto \mathbb{R}, f(x) = \theta^\top x\}$ , we have  $\dim_E(\mathcal{F}, \epsilon) = O(d \log 1/\epsilon)$  since any d linear independent input's estimated value can fully describe a linear mapping function. We also omit the  $\epsilon$  and use  $\dim_E(\mathcal{F})$  when it only has a logarithm dependent term on  $\epsilon$ .

#### 4 Main Results for Contextual Bandits

In this section, we will present our theoretical analysis on the proposed GFUCB algorithm for contextual bandits.

#### 4.1 Assumptions

This section will list the assumptions that we make for our analysis. The main assumption is the existence of a shared feature extraction function from class  $\Phi = \{\phi : \mathcal{C} \times \mathcal{A} \mapsto \mathbb{R}^k\}$  that any task's value function is linear in this  $\phi$ .

**Assumption 1.1 (Shared Space and Representation)** All the tasks share the same context space  $\mathcal C$  and action space  $\mathcal A$ . Also, there exists a shared representation function  $\phi \in \Phi$  and a set of k-dimensional parameters  $\{\theta_i\}_{i=1}^M$  such that each  $f^{(i)}$  has the form  $f^{(i)}(\cdot,\cdot) = \langle \phi(\cdot,\cdot), \theta_i \rangle$ .

Following standard regularization assumptions for bandits [17, 40], we make assumptions on noise distribution and function parameters.

Assumption 1.2 (Conditional Sub-Gaussian Noise) Denote  $\mathcal{H}_{t,i} = \sigma(C_{1,i}, A_{1,i}, \dots, C_{t,i}, A_{t,i})$  to be the  $\sigma$ -field summarizing the history information available before reward  $R_{t,i}$  is observed for every task  $i \in [M]$ . We have  $\eta_{t,i}$  is sampled from a 1-Sub-Gaussian distribution, namely  $\mathbb{E}\left[\exp(\lambda\eta_{t,i}) \mid \mathcal{H}_{t,i}\right] \leq \exp\left(\frac{\lambda^2}{2}\right)$  for  $\forall \lambda \in \mathbb{R}$ 

**Assumption 1.3 (Bounded-Norm Feature and Parameter)** We assume that the parameter  $\theta_i$  and the feature vector for any context-action pair  $(C,A) \in \mathcal{C} \times \mathcal{A}$  is constant bounded for each task  $i \in [M]$ , namely  $\|\boldsymbol{\theta}_i\|_2 \leq \sqrt{k}$  for  $\forall i \in [M]$  and  $\|\phi(C,A)\|_2 \leq 1$  for  $\forall C \in \mathcal{C}, A \in \mathcal{A}$ .

Apart from these assumptions, we add assumption to measure and constrain the complexity of value approximation function class  $\mathcal{F} = \mathcal{L} \circ \Phi$ .

**Assumption 1.4 (Bounded Eluder Dimension).** We assume that function class  $\mathcal{F}$  has bounded Eluder dimension d, which means for any  $\epsilon$ ,  $\dim_E(\mathcal{F}, \epsilon) = \hat{O}(d)$ .

#### 4.2 Algorithm Details

## **Algorithm 1** Generalized Functional UCB Algorithm

- 1: for step  $t: 1 \rightarrow T$  do
- Compute  $\mathcal{F}_t$  according to (\*)
- Receive contexts  $C_{t,i}$  and action sets  $A_{t,i}$ ,  $i \in [M]$ 3:
- $\begin{array}{l} f_t, A_{t,i} = \operatorname{argmax}_{f \in \mathcal{F}_t, \ A_i \in \mathcal{A}_{t,i}} \sum_{i=1}^M f^{(i)}(C_{t,i}, A_i) \\ \text{Play } A_{t,i} \text{ for task i, and get reward } R_{t,i} \text{ for } i \in [M]. \end{array}$
- 6: end for

The details of the algorithm is in Algorithm 1. At each step t, the algorithm first solves the optimization problem below to get the empirically optimal solution  $f_t$  that best predicts the rewards for contextinput pairs seen so far.

$$\hat{f}_t \leftarrow \underset{f \in \mathcal{F}^{\otimes M}}{\operatorname{argmin}} \sum_{i=1}^{M} \sum_{k=1}^{t-1} \left( f^{(i)}(C_{k,i}, A_{k,i}) - R_{k,i} \right)^2$$

Here we abuse the notation of  $\mathcal{F}^{\otimes M}$  as  $\mathcal{F}^{\otimes M} = \{ f = (f^{(1)}, \dots, f^{(M)}) : f^{(i)}(\cdot) = \phi(\cdot)^{\top} w_i \in \mathcal{F} \}$ to denote the M-head prediction version of  $\mathcal{F}$ , parametrized by a shared representation function  $\phi(\cdot)$ and a weight matrix  $W = [w_1, \dots, w_M] \in \mathbb{R}^{k \times M}$ . We use  $f^{(i)}$  to denote the  $i_{th}$  head of function fwhich specially serves for task i.

After obtaining  $\hat{f}_t$ , we maintain a functional confidence set  $\mathcal{F}_t \subseteq \mathcal{F}^{\otimes M}$  for possible value approxi-

$$\mathcal{F}_{t} \stackrel{\text{def}}{=} \left\{ f \in \mathcal{F}^{\otimes M} : \left\| \hat{f}_{t} - f \right\|_{2, E_{t}}^{2} \leq \beta_{t}, |f^{(i)}(\boldsymbol{x})| \leq 1, \forall \boldsymbol{x} \in \mathcal{C} \times \mathcal{A}, i \in [M] \right\}$$
 (\*)

Here, for the sake of simplicity, we use  $\left\|\hat{f}_t - f\right\|_{2,E_t}^2 = \sum_{i=1}^M \sum_{k=1}^{t-1} \left(\hat{f}_t^{(i)}(\boldsymbol{x}_{k,i}) - f^{(i)}(\boldsymbol{x}_{k,i})\right)^2$  to denote the empirical 2-norm of function  $\hat{f}_t - f = \left(\hat{f}_t^{(1)} - f^{(1)}, \dots, \hat{f}_t^{(M)} - f^{(M)}\right)$ . Basically, (\*) contains all the functions in  $\mathcal{F}^{\otimes M}$  whose value estimation difference on all collected contextaction pairs  $x_{k,i} = (C_{k,i}, A_{k,i})$  compared with empirical loss minimizer  $f_t$  does not exceed a preset parameter  $\beta_t$ . We show that with high probability, the real value function  $f_{\theta}$  is always contained in  $\mathcal{F}_t$  when  $\beta_t$  is carefully chosen as  $O(Mk + \log (\mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty})))$ , where  $\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_{\infty})$  is the  $\alpha$ -covering number of function class  $\Phi$  in the sup-norm  $\|\phi\|_{\infty} = \max_{x \in \mathcal{S} \times \mathcal{A}} \|\phi(x)\|_2$  and  $\alpha$  is set to be a small number as  $\frac{1}{kMT}$  (see detailed definition and proof in Lemma 1).

For the action choice, our algorithm follows OFUL, which estimates each action value with the most optimistic function value in our confidence set  $\mathcal{F}_t$ , and chooses the action whose optimistic value estimation is the highest. In the multitask setting, we choose one action from each task to form an action tuple  $(A_1, A_2, \dots, A_M)$  such that the summation of the optimistic value estimation  $\sum_{i=1}^{M} f^{(i)}(C_{t,i}, A_i)$  is maximized by some function  $f \in \mathcal{F}_t$ .

Intractability. Some may have concerns on the intractability of building the confidence set (\*) and solving the optimization problem to get  $f_t$ ,  $f_t$ ,  $A_{t,i}$ . The solution comes as two folds. From the theoretical perspective, since the focus of problem is sample complexity rather than computational complexity, a computational oracle can simply be assumed to give the solution of the optimization. This is the common practice for theoretical works [20, 35, 1, 19] in order to focus on the sample complexity analysis. From empirical perspective, there are great chances to optimize it with gradient methods. For example, solving  $\hat{f}_t$  is a standard empirical risk minimization problem, and can be effectively solved with gradient methods [13]. As for  $f_t$  and  $A_{t,i}$ , note that it is not necessary to explicitly build the confidence set  $\mathcal{F}_t$  by listing all the candidates. The approximation algorithm just need to search within the confidence set via gradient method to optimize objective  $\sum_{i=1}^M f^{(i)}(C_{t,i}, A_i)$ . The start point is  $\hat{f}_t$ , and the algorithm knows that it approaches the border of  $\mathcal{F}_t$  when  $\|\hat{f}_t - f\|_{2,E_t}^2$  approaches  $\beta_t$ . The details of implementation are in section 6.

**Mechanism.** GFUCB algorithm solves the exploration problem in an implicit way. For a context-action pair  $\boldsymbol{x} = (C, A)$  in task i which has not been fully understood and explored yet, the possible value estimation  $f^{(i)}(\boldsymbol{x})$  will vary in large range with regard to constraint  $\|f - \hat{f}_t\|_{2,E_t}^2 \leq \beta_t$ . This is because within  $\mathcal{F}_t$  there are many possible function value on this  $\boldsymbol{x}$  while agreeing on all past context-action pairs' value. Therefore, the optimistic value  $f^{(i)}(\boldsymbol{x})$  will become high by getting a significant implicit bonus, encouraging the agent to try such action A under context C, which achieves natural exploration.

The reduction of sample complexity is achieved through joint training for function  $\phi$ . If we solve these tasks independently, the confidence set width  $\beta_t$  is at scale  $M\log\left(\mathcal{N}(\Phi,\alpha,\|\cdot\|_\infty)\right)$  because it needs to cover M representation function space respectively. By involving  $\phi$  in the prediction for all tasks, our algorithm reduces the size of confidence set by M times, since now the samples from all the tasks can contribute to learn the representation  $\phi$ . Usually  $\log\left(\mathcal{N}(\Phi,\alpha,\|\cdot\|_\infty)\right)$  is much greater than k and M, hence our confidence set shrinks at a much faster speed. This explains how GFUCB achieves lower regret, since the sub-optimality at each step t is proportional to the confidence set width  $\beta_t$  when real value function  $f_\theta \in \mathcal{F}_t$ .

## 4.3 Regret Bound

Based on the assumptions above, we have the regret guarantee as below.

**Theorem 1.** Based on assumption 1.1 to 1.4, denote the cumulative regret in T steps as  $\operatorname{Reg}(T)$ , with probability at least  $1 - \delta$  we have  $\operatorname{Reg}(T) = \tilde{O}\left(\sqrt{MdT(Mk + \log \mathcal{N}(\Phi, \alpha_T, \|\cdot\|_{\infty}))}\right)$ .

Here,  $d := \dim_E(\mathcal{F}, \alpha_T)$  is the Eluder dimension for value approximation function class  $\mathcal{F} = \mathcal{L} \circ \Phi$ , and  $\alpha_T$  is discretization scale which only appears in logarithm term thus omitted. The detailed proof is left in appendix.

To the best of knowledge, this is the first regret bound for general function class representation learning in contextual bandits. To get a sense of its sharpness, note that when  $\Phi$  is specialized as linear function class as  $\Phi = \{\phi(x) = \boldsymbol{B}\boldsymbol{x}, \boldsymbol{B} \in \mathbb{R}^{k \times d}\}$ , we have  $\log \mathcal{N}(\Phi, \alpha_T, \|\cdot\|_{\infty}) = \tilde{O}(dk)$  and  $\dim_E(\mathcal{F}) = d$ , then our bound is reduced to  $\tilde{O}(M\sqrt{dTk} + d\sqrt{MTk})$ , which is the same optimal as the current best provable regret bound for linear representation class bandits in [17].

## 5 Main Results for MDP

## 5.1 Assumptions

For multitask Linear MDP setting, we adopt Assumption 3 from [17] which generalizes the inherent Bellman error [44] to multitask setting.

Assumption 2.1 (Low IBE for multitask) Define multi-task IBE is defined as

$$\mathcal{I}_{h}^{\text{mul}} \overset{\text{def}}{=} \sup_{\left\{Q_{h+1}^{(i)}\right\}_{i=1}^{M} \in \mathcal{Q}_{h+1}} \inf_{\left\{Q_{h}^{(i)}\right\}_{i=1}^{M} \in \mathcal{Q}_{h}} \sup_{s \in \mathcal{S}, a \in \mathcal{A}, i \in [M]} \left| \left(Q_{h}^{(i)} - \mathcal{T}_{h}^{(i)} \left(Q_{h+1}^{(i)}\right)\right)(s, a) \right|.$$

We have  $\mathcal{I} \stackrel{\text{def}}{=} \sup_h \mathcal{I}_h^{\text{mul}}$  is small for all  $\mathcal{Q}_h$ ,  $h \in [H]$ .

Assumption 2.1 generalize low IBE to multitask setting. It assumes that for every task  $i \in [M]$ , its Q-value function space is always close under Bellman operator.

## Assumption 2.2 (Parameter Regularization) We assume that

- $\|\phi(s,a)\| \le 1$ ,  $0 \le Q_h^{\pi}(s,a) \le 1$  for  $\forall (s,a) \in \mathcal{S} \times \mathcal{A}$ ,  $h \in [H], \forall \pi$ .
- There exists a constant D such that for any  $h \in [H]$  and  $\theta_h^{(i)}$ , it holds that  $\|\theta_h^{(i)}\|_2 \leq D$ .
- For any fixed  $\left\{Q_{h+1}^{(i)}\right\}_{i=1}^{M} \in \mathcal{Q}_{h+1}$ , the random noise  $z_h^{(i)} \stackrel{\text{def}}{=} R_h^{(i)}(s,a) + \max_a Q_{h+1}^{(i)}(s',a) \mathcal{T}_h^{(i)}\left(Q_{h+1}^{(i)}\right)(s,a)$  is bounded in [-1,1] and is always independent to all other random variables for  $\forall (s,a) \in \mathcal{S} \times \mathcal{A}, h \in [H], i \in [M]$ .

These assumptions are widely adopted in linear MDP analytical works [44, 17, 26], which regularizes the parameter, feature, and noise scale. Again we add bounded Eluder dimension constraint for the Q-value estimation class.

**Assumption 2.3 (Bounded Eluder Dimension).** We assume that function class  $Q_h$  has bounded Eluder dimension d for any  $h \in [H]$ .

## 5.2 Algorithm Details

## Algorithm 2 multitask Linear MDP Algorithm

```
1: for episode t: 1 \to T do

2: Q_{H+1}^{(i)} = 0, i \in [M]

3: for h: H \to 1 do

4: \hat{\phi}_{h,t}, \hat{\theta}_{h,t}^{(i)} \leftarrow \text{solving } (1)

5: Q_{h}^{(i)}(\cdot, \cdot) = \hat{\phi}_{h,t}(\cdot, \cdot)^{\top} \hat{\theta}_{h,t}^{(i)}, V_{h}^{(i)}(\cdot) = \max_{a} Q_{h}^{(i)}(\cdot, a)

6: end for

7: for h: 1 \to H do

8: Compute \mathcal{F}_{h,t} according to Lemma 4

9: Receive states \left\{s_{h,t}^{(i)}\right\}_{i=1}^{M}, \tilde{f}_{h,t}, a_{h,t}^{(i)} = \operatorname{argmax}_{f \in \mathcal{F}_{h,t}, a^{(i)} \in \mathcal{A}} \sum_{i=1}^{M} f^{(i)}\left(s_{h,t}^{(i)}, a^{(i)}\right)

10: Play a_{h,t}^{(i)} and get reward R_{h,t}^{(i)} for task i \in [M].

11: end for

12: end for
```

The algorithm for multitask linear MDP is similar to contextual bandits as above. The optimization problem in line 4 of Algorithm 2 is finding the empirically best solution for Q-value estimation at level h in episode t as below

$$\hat{\phi}_{h,t}, \hat{\boldsymbol{\Theta}}_{h,t} \leftarrow \underset{\boldsymbol{\phi} \in \Phi, \boldsymbol{\Theta} = [\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(M)}]}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\Theta})$$

$$s.t. \quad \|\boldsymbol{\theta}^{(i)}\| \leq D, \forall i \in [M]$$

$$0 \leq \boldsymbol{\phi}(s, a)^{\top} \boldsymbol{\theta}_{i} \leq 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, i \in [M],$$

$$(1)$$

where  $\mathcal{L}(\phi, \boldsymbol{\Theta})$  is the empirical loss function defined as

$$\sum_{i=1}^{M} \sum_{j=1}^{t-1} \left( \phi \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right)^{\mathsf{T}} \boldsymbol{\theta}^{(i)} - R_{h,j}^{(i)} - V_{h+1}^{(i)} \left( s_{h+1,j}^{(i)} \right) \right)^{2}.$$

The framework of our work resembles LSVI [21] and [26] which learns the Q-value estimation in a reverse order, at each level h, the algorithm uses just-learned value estimation function  $V_{h+1}$  to build the regression target value as  $R_{h,j}^{(i)} + V_{h+1}^{(i)} \left(s_{h+1,j}^{(i)}\right)$  and find empirically best estimation  $\hat{f}_{h,t}^{(i)} = \hat{\phi}_{h,t}^{\top} \hat{\theta}_{h,t}^{(i)}$  for each task  $i \in [M]$ . The optimistic value estimation of each action is again searched within confidence set  $\mathcal{F}_{h,t}$  which centered at  $\hat{f}_{h,t}$  and shrinks as the constraint  $\|f - \hat{f}_{h,t}\|_{2,E_t}^2 \leq \beta_t$  becomes increasingly tighter. Note that the contextual bandit problem can be regarded as a 1-horizon MDP problem without transition dynamics, and our framework at each level h is indeed a copy of procedures in Algorithm 1.

#### 5.3 Regret Bound

Based on assumptions 2.1 to 2.3, we prove that our algorithm enjoys a regret bound guaranteed by the following theorem. Detailed proof is left in appendix.

**Theorem 2.** Based on assumption 2.1 to 2.3, denote the cumulative regret in T episodes as Reg(T), we have the following regret bound for Reg(T) holds with probability at least  $1 - \delta$  for Algorithm 2

$$\tilde{O}\left(MH\sqrt{Tdk} + H\sqrt{MTd\log\mathcal{N}(\Phi,\alpha)} + MHT\mathcal{I}\sqrt{d}\right),$$

where  $\alpha$  is discretization scale smaller than  $\frac{1}{kMT}$ .

**Remark.** Compared with naively executing single task general value function approximation algorithm [39] for M tasks, whose regret bound is  $\tilde{O}(MHd\sqrt{T\log\mathcal{N}(\Phi)})$ , to achieve same average regret, our algorithm outperforms this naive algorithm with a boost of sample efficiency by  $\tilde{O}(Md)$ . This benefit mainly attributes to learning in function space  $\mathcal{F}^{\otimes M} = \mathcal{L}^M \circ \Phi$  instead of  $\mathcal{F}^M = (\mathcal{L} \circ \Phi)^M$ , the former is more compact and requires much less samples to learn.

## 6 Experiments

To validate our theoretical findings, we conduct experiments on a non-linear neural network bandits. Note that it is a proof-of-concept experiment. Our main purpose is to realize the GFUCB algorithm and check its efficacy but *not* to beat sophisticated real-world algorithms. The point to demonstrate is that sample efficiency of GFUCB is scalable to the number of tasks and better than naive exploration.

#### 6.1 Task Design

To test the efficacy of our algorithm, we use the MNIST dataset [10] to build a bandit problem that involves non-linear value approximation. The reward function of the bandit environment maps the same digit into the same base reward  $r_b$ , which ranges from 0 to 1, plus a noise  $\eta_h$  sampled from a zero-mean Gaussian with a standard deviation of 0.01. At every round, each task will present the agent a context C consists of K different digit images and ask the agent to take action as an integer  $i \in [K]$  meaning which image to choose, then return the reward according to the agent's choice.

For the multitask setting, we construct M different tasks using different digit-to-reward mappings  $\sigma_i:\{0,\ldots,9\}\mapsto [0,1], i\in [M],$  where  $\sigma_i(k)$  will give a unique reward for all images of digit k in task i. Different tasks have different reward mapping function  $\sigma_i(\cdot)$ . By designing the environment this way, it requires to learn a common representation  $\phi$  to recognize digits for different tasks.

#### **6.2** Implementation Details

We use a simple CNN as our feature extraction function  $\phi$ , which takes a digit image as input and outputs a 10-dimensional normalized vector as representation. It consists of two 3x3 convolution layers and two fully-connected layers, followed by ReLU activation and a normalization procedure.

The biggest challenge for implementation is how to solve a complex optimization problem in general functional space. In principle, finding parameters for a neural network to achieve the (near) minimal empirical error is an NP-Hard problem. To solve this issue, we use a gradient-based method to approximately find a local-optimal solution. For finding the empirically best  $\hat{f}_t$ , we use Adam with lr=1e-3 to train for sufficiently long steps; in our setting, it is set to be 200 epochs at every step t, to ensure that the training loss is sufficiently low.

The next major challenge is estimating the optimistic value for each action within the abstract function set  $\mathcal{F}_t$ . To tackle this problem, we enumerate all possible action tuples  $\{A_i\}_{i=1}^M$  and then solve the equivalent optimization below to compute its optimistic estimated value

$$\max_{f \in \mathcal{F}_t} \sum_{i=1}^{M} f^{(i)}(C_{t,i}, A_i) \quad s.t. \quad \left\| f - \hat{f}_t \right\|_{2, E_t}^2 \le \beta_t.$$

Still, this is a complicated optimization problem within an abstract function set. Inspired by the Lagrangian operator, we transform it into an unconstrained optimization problem minimizing loss

function  $\ell(f) = -\sum_{i=1}^M f^{(i)}(A_i) + \lambda \cdot \max(0, \|\hat{f}_t - f\|_{2, E_t}^2 - B_t)$ , where  $\lambda$  is a hyperparameter to be determined, in our algorithm we set it to be  $\lambda = 30$  by empirical search. Also  $B_t = a \log(b \cdot t + c)$  is an approximation for  $\beta_t$  since  $\beta_t$  includes  $\mathcal{N}(\Phi, \alpha)$  which is intractable to be exactly computed, we found (a, b, c) = (0.4, 0.5, 2) to be a good parameter of UCB in single task. We use SGD with a small learning rate (5e - 4) to finetune the model  $\hat{f}_t$  for 200 iterations to optimize  $\ell(f)$ .

The basic intuition is that, through optimizing  $\ell(f)$ , the algorithm will try to maximize function value  $\sum_{i=1}^M f^{(i)}(A_i)$ . And as long as f satisfies  $\|\hat{f}_t - f\|_{2,E_t}^2 \leq B_t$ , such constraint will not appear in the loss term, thus has no effect on optimization. When f comes to the border of  $\mathcal{F}_t$ , where  $\|\hat{f}_t - f\|_{2,E_t}^2$  approaches  $B_t$ , the second term adds regularization term to the loss as punishment, preserving  $\|\hat{f}_t - f\|_{2,E_t}^2$  at a near-constant level around  $B_t$ . So we can approximately simulate the optimistic value estimating procedure via searching in the neighborhood of  $\hat{f}_t$ .

#### 6.3 Connection to Algorithm 1

The main difference between our practical version algorithm and the theoretical one is that we did not list out all the functions in the whole confidence set  $\mathcal{F}_t$  explicitly, but just use gradient-based method to implicitly search within a very small fraction of  $\mathcal{F}_t$  with heuristics. Getting a candidate within the confidence set is much easier and tractable than rigorously exhausting all functions in  $\mathcal{F}_t$  to optimize. We can start from the parameter of  $\hat{f}_t$  and use gradient method to approximately find  $f_t$  and  $A_{t,i}$ .

Another difference is we do not rigorous compute  $\beta_t$  which involves  $\mathcal{N}(\Phi)$ , but directly determine a parametrized function form. Rigorously speaking, our tuned value of  $\beta_t$  is much smaller than the theoretical guaranteed ones, so all the candidate functions that we search along the trajectory of gradient method still satisfy the theoretical requirement (but it may omit many other potential candidates). Therefore, our practical version algorithm should be regarded as an inaccurate approximation to the theoretical algorithm. Moreover, it also plays a role as regularization to enable the convergence of  $\mathcal{F}_t$  since we only consider regular ones in the neighborhood of  $\hat{f}_t$ .

#### 6.4 Results

We test the performance of our algorithm against a naive eps-greedy baseline that solves each task independently by training the same CNN value prediction module. We show our results with number of tasks M=1,5,10 in Figure 1. Firstly, we randomly generate 10 different digit-value mapping functions  $\sigma_i(\cdot), i=1,\ldots,10$ . The total 10 tasks are divided into 10/M groups; each group forms a M-task problem and is solved by an individual copy of some algorithm. At each step t, the cumulative regret from all 10 tasks is averaged to estimate the method's performance. Our result in Figure 1 verified that the multitask training does accelerate learning, which empirically validates our theoretical analysis. The multitask training utilizes the samples from all M tasks to jointly learn a good represen-

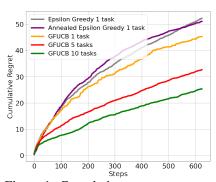


Figure 1: Cumulative regret over steps for M=1,5,10.

tation  $\phi$ , which significantly accelerates the learning procedure of the CNN backbone. Also, the improvement in GFUCB algorithm's performance with M=1 validates the effect of our finetune procedure for getting a bonus. Detailed dissection and discussion are left in appendix.

#### 7 Conclusion

In this work, we extend the analysis of the benefit of multitask representation learning from linear representation class to general function class. We propose a straightforward algorithm that can utilize samples from all the tasks to jointly train a representation function, which is demonstrated theoretically and empirically to accelerate the sample efficiency and outperform naively single-task learning. Also, we extend the analysis to the MDP setting and show that the benefit of multitask representation learning is similar. Furthermore, our experimental result reveals that our proposed algorithm is also effective in practice even for highly non-linear neural network representations.

## **Acknowledgments and Disclosure of Funding**

This work is supported in part by the National Science and Technology Major Project of the Ministry of Science and Technology of China under Grants 2018AAA0101604, the National Natural Science Foundation of China under Grants 62022048 and the State Key Lab of Autonomous Intelligent Unmanned Systems.

#### References

- [1] Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pages 1638–1646. PMLR, 2014.
- [2] Rie Kubota Ando and Tong Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(Nov):1817–1853, 2005.
- [3] Sanjeev Arora, Simon S Du, Sham Kakade, Yuping Luo, and Nikunj Saunshi. Provable representation learning for imitation learning via bi-level optimization. *arXiv preprint arXiv:2002.10544*, 2020.
- [4] Jonathan Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 12:149–198, 2000.
- [5] Shai Ben-David and Reba Schuller. Exploiting task relatedness for multiple task learning. In *Learning Theory and Kernel Machines*, pages 567–580. Springer, 2003.
- [6] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv* preprint arXiv:1912.06680, 2019.
- [7] Rich Caruana. Multitask learning. Machine learning, 28(1):41–75, 1997.
- [8] Giovanni Cavallanti, Nicolo Cesa-Bianchi, and Claudio Gentile. Linear algorithms for online multitask classification. *Journal of Machine Learning Research*, 11(Oct):2901–2934, 2010.
- [9] Yuan Cheng, Songtao Feng, Jing Yang, Hong Zhang, and Yingbin Liang. Provable benefit of multitask representation learning in reinforcement learning. *arXiv preprint arXiv:2206.05900*, 2022.
- [10] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [11] Carlo D'Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Sharing knowledge in multi-task deep reinforcement learning. In *International Conference on Learning Representations*, 2020.
- [12] 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.
- [13] 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*, pages 1675–1685. PMLR, 2019.
- [14] 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.
- [15] Simon S Du, Jayanth Koushik, Aarti Singh, and Barnabás Póczos. Hypothesis transfer learning via transformation functions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 574–584, 2017.
- [16] Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado van Hasselt. Multi-task deep reinforcement learning with popart. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3796–3803, 2019.
- [17] Jiachen Hu, Xiaoyu Chen, Chi Jin, Lihong Li, and Liwei Wang. Near-optimal representation learning for linear bandits and linear rl. In *International Conference on Machine Learning*, pages 4349–4358. PMLR, 2021.

- [18] Yao Hu, Debing Zhang, Jieping Ye, Xuelong Li, and Xiaofei He. Fast and accurate matrix completion via truncated nuclear norm regularization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(9):2117–2130, 2013.
- [19] 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-Volume 70*, pages 1704–1713. JMLR. org, 2017.
- [20] Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418, 2021.
- [21] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv* preprint arXiv:1907.05388, 2019.
- [22] Alessandro Lazaric and Marcello Restelli. Transfer from multiple mdps. In *Advances in Neural Information Processing Systems*, pages 1746–1754, 2011.
- [23] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [24] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [25] Lydia T Liu, Urun Dogan, and Katja Hofmann. Decoding multitask dqn in the world of minecraft. In The 13th European Workshop on Reinforcement Learning (EWRL) 2016, 2016.
- [26] Rui Lu, Gao Huang, and Simon S Du. On the power of multitask representation learning in linear mdp. *arXiv preprint arXiv:2106.08053*, 2021.
- [27] Andreas Maurer. Bounds for linear multi-task learning. *Journal of Machine Learning Research*, 7(Jan):117–139, 2006.
- [28] Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- [29] 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.
- [30] Matteo Papini, Andrea Tirinzoni, Aldo Pacchiano, Marcello Restelli, Alessandro Lazaric, and Matteo Pirotta. Reinforcement learning in linear mdps: Constant regret and representation selection. *Advances in Neural Information Processing Systems*, 34:16371–16383, 2021.
- [31] Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Actor-mimic: Deep multitask and transfer reinforcement learning. *arXiv preprint arXiv:1511.06342*, 2015.
- [32] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *NIPS*, pages 2256–2264. Citeseer, 2013.
- [33] Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy distillation. *arXiv preprint arXiv:1511.06295*, 2015.
- [34] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [35] Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based reinforcement learning in contextual decision processes. arXiv preprint arXiv:1811.08540, 2018.
- [36] Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(Jul):1633–1685, 2009.
- [37] Yee Teh, Victor Bapst, Wojciech M Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas Heess, and Razvan Pascanu. Distral: Robust multitask reinforcement learning. In Advances in Neural Information Processing Systems, pages 4496–4506, 2017.
- [38] Nilesh Tripuraneni, Chi Jin, and Michael I Jordan. Provable meta-learning of linear representations. *arXiv preprint arXiv:2002.11684*, 2020.

- [39] Ruosong Wang, Ruslan Salakhutdinov, and Lin F Yang. Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *arXiv* preprint arXiv:2005.10804, 2020.
- [40] Jiaqi Yang, Wei Hu, Jason D. Lee, and Simon Shaolei Du. Impact of representation learning in linear bandits. In *International Conference on Learning Representations*, 2021.
- [41] Jiaqi Yang, Qi Lei, Jason D Lee, and Simon S Du. Nearly minimax algorithms for linear bandits with shared representation. *arXiv preprint arXiv:2203.15664*, 2022.
- [42] Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.
- [43] Lin F Yang, Chengzhuo Ni, and Mengdi Wang. Learning to control in metric space with optimal regret. *arXiv preprint arXiv:1905.01576*, 2019.
- [44] Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pages 10978–10989. PMLR, 2020.

## A Bandit Regret Bound Analysis

## A.1 Algorithm Procedure

At each round  $s \in [t]$ , after performing a list of actions  $\{A_{s,i}\}_{i=1}^M$  with respect to corresponding context vectors  $\{C_{s,i}\}_{i=1}^M$ , the agent receives a list of rewards  $y_{s,i}$  associated with input  $\boldsymbol{x}_{s,i} = (C_{s,i}, A_{s,i})$  for  $i \in [M]$ . Note that we will use  $f(C_t, A_t)$  or  $f(\boldsymbol{x}_t)$  where  $\boldsymbol{x}_t = (C_t, A_t)$  in different contexts. The algorithm first solves the following regression problem to obtain the empirical minimizer function  $\hat{f}_t(\cdot) = \hat{\phi}_t(\cdot)^{\top} \widehat{\boldsymbol{W}}_t$  based on samples collected.

$$\widehat{\phi}_{t}, \widehat{\boldsymbol{W}}_{t} = \operatorname*{argmin}_{\phi \in \Phi, \boldsymbol{W} = [\boldsymbol{w}_{1}, \dots, M]} \sum_{i=1}^{M} \|\boldsymbol{y}_{t-1, i} - \phi(\boldsymbol{X}_{t-1, i})^{\top} \boldsymbol{w}_{i}\|_{2}^{2}$$

$$s.t. \quad |\phi(\boldsymbol{x})^{\top} \boldsymbol{w}_{i}| \leq 1, \quad \forall i \in [M], \boldsymbol{x} \in \mathcal{C} \times \mathcal{A}.$$

Here,  $\boldsymbol{X}_{t-1,i} = [\boldsymbol{x}_{1,i}, \boldsymbol{x}_{2,i}, \dots, \boldsymbol{x}_{t-1,i}]$  is the selected context-action pair for task i in the first t-1 rounds, and  $\boldsymbol{y}_{t-1,i} = [R_{1,i}, R_{2,i}, \dots, R_{t-1,i}]^{\top} \in \mathbb{R}^{t-1}$  stacks all the received reward into a vector accordingly. We use  $\phi(\boldsymbol{X})$  to compactly represent feeding each column  $\boldsymbol{x}_i$  of  $\boldsymbol{X}$  into  $\phi(\cdot)$  and get concatenated output as  $[\phi(\boldsymbol{x}_1), \phi(\boldsymbol{x}_2), \dots, \phi(\boldsymbol{x}_{t-1})]$ .

After obtaining the best empirical estimator function  $\hat{f}_t^{(i)}(\cdot) = \hat{\phi}_t(\cdot)^\top \hat{w}_{t,i}$  at round  $t \in [T]$  for each  $i \in [M]$ , we maintain a function confidence set  $\mathcal{F}_t \subseteq \mathcal{F}^{\otimes M}$  for representation function and parameters.

$$\mathcal{F}_{t} \stackrel{\text{def}}{=} \left\{ f \in \mathcal{F}^{\otimes M} : \left\| \hat{f}_{t} - f \right\|_{2, E_{t}}^{2} \leq \beta_{t}, |f^{(i)}(\boldsymbol{x})| \leq 1, \forall \boldsymbol{x} \in \mathcal{C} \times \mathcal{A}, i \in [M] \right\}$$

$$(*)$$

Here we abuse the notation of  $\mathcal{F}^{\otimes M}$  as  $\mathcal{F}^{\otimes M} = \left\{ f = \left( f^{(1)}, \dots, f^{(M)} \right) : f^i(\cdot) = \phi(\cdot)^\top \boldsymbol{w}_i \in \mathcal{F} \right\}$  to denote the M-head prediction version of  $\mathcal{F}$ , parametrized by a shared representation function  $\phi(\cdot)$  and a weight matrix  $\boldsymbol{W} = [\boldsymbol{w}_1, \dots, \boldsymbol{w}_M] \in \mathbb{R}^{k \times M}$ . We use  $f^{(i)}$  to denote the  $i_{th}$  head of function f. For the sake of simplicity, we use

$$\left\| \hat{f}_t - f \right\|_{2, E_t}^2 = \sum_{i=1}^M \sum_{s=1}^{t-1} \left( \hat{f}_t^{(i)}(\boldsymbol{x}_{s,i}) - f^{(i)}(\boldsymbol{x}_{s,i}) \right)^2$$

to denote the empirical 2-norm of function  $\hat{f}_t - f = (\hat{f}_t^{(1)} - f^{(1)}, \dots, \hat{f}_t^{(M)} - f^{(M)})$ . Another important hyperparameter for our algorithm is the confidence set width term  $\beta_t$ , which is a function of representation function class  $\Phi$ , probability  $\delta$  and discretization scale parameter  $\alpha$ .

$$\beta_t(\Phi, \alpha, \delta) = 12Mk + 12\log\left(\mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty})/\delta\right) + 8\alpha\sqrt{Mtk(Mt + \log(2Mt^2/\delta))}$$

here  $\mathcal{N}(\mathcal{F}, \alpha, \|\cdot\|_{\infty})$  is the  $\alpha$ -covering number of function class  $\Phi$  in the sup-norm  $\|\phi\|_{\infty} = \max_{\boldsymbol{x} \in \mathcal{S} \times \mathcal{A}} \|\phi(\boldsymbol{x})\|_2$  (see detailed definition in Lemma 1) and  $\alpha$  can be set to be some small scale number, like  $\frac{1}{kMT}$ .

#### A.2 Main Proof sketch

In this section we will give a theoretical guarantee for the performance of our algorithm. Before diving into details, we first explain the overall idea and structure of our proof. First, we decompose the regret into the summation of confidence set width at different rounds plus a small term which accounts for the possibility that confidence function set  $\mathcal{F}_t$  fails to contain ground truth function  $f_{\theta}$ .

**Lemma 0.** Fix any sequence of confidence set  $\{\mathcal{F}_t, t \in \mathbb{N}\}$  which is measurable with respect to history  $\mathcal{H}_t$ , denote the induced policy by Algorithm 1 as  $\pi = \{\pi_i\}_{i=1}^M$  where each  $\pi_i : \mathcal{C} \mapsto \mathcal{A}, i \in [M]$  is for task i, then for any  $T \in \mathbb{N}$  we have

$$\operatorname{Regret}(T) := \sum_{i=1}^{M} \sum_{t=1}^{T} f_{\theta}^{(i)} \left( \boldsymbol{x}_{t,i}^{\star} \right) - f_{\theta}^{(i)} (\boldsymbol{x}_{t,i}) \leq \sum_{t=1}^{T} \left[ w_{\mathcal{F}_{t}}(\boldsymbol{X}_{t}) + C \cdot \mathbb{I}(f_{\theta} \notin \mathcal{F}_{t}) \right]$$

where  $x_{t,i} = (C_{t,i}, \pi_i(C_{t,i}))$  is the context-action pair that actually happened.  $A_{t,i}^\star = \arg\max_A f_{\theta}^{(i)}(C_{t,i}, A)$  is the optimal action for each task  $i \in [M]$  at round  $t \in [T]$ , and  $x_{t,i}^\star = (C_{t,i}, A_{t,i}^\star)$  is the corresponding optimal context-action pair, C is a universal large enough constant. We use  $X_t = [x_{t,1}, \dots, x_{t,M}]$  to stack  $x_{t,i}$  into a matrix, similar for  $X_t^\star = [x_{t,1}^\star, \dots, x_{t,M}^\star]$ . The confidence set width  $w_{\mathcal{F}_t}(X_t)$  is defined by

$$w_{\mathcal{F}_t}(\boldsymbol{X}_t) := \sup_{\overline{f}, f \in \mathcal{F}_t} \sum_{i=1}^M \left[ \overline{f}^{(i)}(\boldsymbol{x}_{t,i}) - \underline{f}^{(i)}(\boldsymbol{x}_{t,i}) \right].$$

Essentially, it measures the largest total difference of value estimation among all the functions in  $f \in \mathcal{F}_t$  for the fixed inputs  $x_{t,i}$  where  $i \in [M]$ . Apart from the constant term accounting for the case that  $\mathcal{F}_t$  fails to contain  $f_{\theta}$ , which we will prove happen with small probability, this regret is then bounded by the sum of width over time step t.

Next, we will show that our construction of confidence set  $\mathcal{F}_t$  makes all of them contain real value function with high probability.

**Lemma 1.** For all  $\delta \in (0,1)$  and  $\alpha > 0$ , if  $\mathcal{F}_t$  is defined by  $\mathcal{F}_t = \{f \in \mathcal{F}^{\otimes M} : ||f - \hat{f}||_{2,E_t} \leq \sqrt{\beta_t(\Phi, \delta, \alpha)}\}$  for all  $t \in \mathbb{N}$ , where  $\hat{f}$  is the solution to the empirical error minimization. Denote the ground truth value function as  $f_{\theta}(\cdot)$ , then we have

$$\mathbb{P}\left(f_{\theta} \in \bigcap_{t=1}^{T} \mathcal{F}_{t}\right) \ge 1 - 2\delta.$$

After that, we prove that

#### Lemma 2.

$$\sum_{t=1}^{T} \mathbb{I}\left(w_{\mathcal{F}_t}(\boldsymbol{X}_t) > \epsilon\right) \le \left(\frac{4M\beta_T}{\epsilon^2} + 1\right) \dim_E(\mathcal{F}, \epsilon)$$

Then plug it into lemma 0, we get our main result for the regret bound as

$$\operatorname{Reg}(\pi, T) \le \frac{1}{T} + \min \left\{ \dim_E(\mathcal{F}, \alpha_T), T \right\} + 4\sqrt{M \dim_E(\mathcal{F}, \alpha_T)\beta_T T}$$
 (1)

Usually  $\alpha_T$  is set to be a small number like  $\frac{1}{kMT}$ , or the minimizer for  $\beta_T(\Phi, \alpha, \delta)$ . We know that  $\dim_E(\mathcal{F}, \alpha_T)$  is a poly-logarithmic function of T, which means the final regret bound is dominant by term  $\sqrt{M\dim_E(\mathcal{F}, \alpha_T)\beta_T T}$  when  $T \to \infty$ . This further becomes

$$\sqrt{MT\left(Mk + \log\left(\mathcal{N}(\Phi, (kMT)^{-1}, \|\cdot\|_{\infty})\right)\right) \dim_{E}(\mathcal{F}, (kMT)^{-1})} \tag{2}$$

For example, if  $\Phi$  is specialized as linear function class parametrized by matrix  $\Theta \in \mathbb{R}^{d \times k}$ , then  $\log \left( \mathcal{N}(\Phi, (kMT)^{-1}, \|\cdot\|_{\infty}) \right) = O(kd \log(kMT))$  and  $\dim_E(\mathcal{F}, (kMT)^{-1}) = O(d \log(kMT))$ , hence the regret bound becomes

$$O(\sqrt{MT(Mk+kd)d}\log(kMT)) = \tilde{O}(M\sqrt{kdT} + d\sqrt{MkT})$$

which reduces to result in [17] by a poly-logarithm factor.

#### A.3 Detailed Proof

Proof of Lemma 0. Define the upper and lower bounds  $U_t(\boldsymbol{X}_t) = \sup \left\{ \sum_{i=1}^M f^{(i)}(\boldsymbol{x}_{t,i}) : f \in \mathcal{F}_t \right\}$  and  $L_t(\boldsymbol{X}_t) = \inf \left\{ \sum_{i=1}^M f^{(i)}(\boldsymbol{x}_{t,i}) : f \in \mathcal{F}_t \right\}$ .

If  $f_{\theta} \notin \mathcal{F}_t$ , then the error will be bounded by a large constant C since all f(x) is constant bounded. Otherwise  $f_{\theta} \in \mathcal{F}_t$ , we have

$$L_t(\boldsymbol{X}_t) \leq \sum_{i=1}^M f_{\theta}^{(i)}(\boldsymbol{x}_{t,i}) \leq U_t(\boldsymbol{X}_t)$$

$$\sum_{i=1}^{M} f_{\theta}^{(i)}(\boldsymbol{x}_{t,i}^{\star}) \leq U_{t}(\boldsymbol{X}_{t}^{\star})$$

where  $X_t$  and  $X_t^*$  is defined in lemma 0. Also, by the optimality of  $X_t$  with respect to  $\mathcal{F}_t$ , we know  $U_t(X_t^*) \leq U_t(X_t)$ , therefore

$$\sum_{i=1}^{M} \left[ f_{\theta}^{(i)}(\boldsymbol{x}_{t,i}^{\star}) - f_{\theta}^{(i)}(\boldsymbol{x}_{t,i}) \right] \leq C \cdot \mathbb{I}(f_{\theta} \notin \mathcal{F}_{t}) + \left[ U_{t}(\boldsymbol{X}_{t}^{\star}) - L_{t}(\boldsymbol{X}_{t}) \right]$$

$$= C \cdot \mathbb{I}(f_{\theta} \notin \mathcal{F}_{t}) + \sum_{i=1}^{M} \left[ U_{t}(\boldsymbol{X}_{t}^{\star}) - U_{t}(\boldsymbol{X}_{t}) + U_{t}(\boldsymbol{X}_{t}) - L_{t}(\boldsymbol{X}_{t}) \right]$$

$$\leq C \cdot \mathbb{I}(f_{\theta} \notin \mathcal{F}_{t}) + \sum_{i=1}^{M} \left[ U_{t}(\boldsymbol{X}_{t}) - L_{t}(\boldsymbol{X}_{t}) \right]$$

$$= C \cdot \mathbb{I}(f_{\theta} \notin \mathcal{F}_{t}) + w_{\mathcal{F}_{t}}(\boldsymbol{X}_{t})$$

Take summation over  $t \in [T]$  and complete the proof.

**Lemma 1.** For all  $\delta \in (0,1)$  and  $\alpha > 0$ , if  $\mathcal{F}_t$  is defined by  $\mathcal{F}_t = \left\{ f \in \mathcal{F}^{\otimes M} : \|f - \hat{f}\|_{2,E_t} \leq \sqrt{\beta_t(\Phi,\delta,\alpha)} \right\}$  for all  $t \in \mathbb{N}$ , where  $\hat{f}$  is the solution to the empirical error minimization. Denote the ground truth value function as  $f_{\theta}$ , then we have

$$\mathbb{P}\left(f_{\theta} \in \bigcap_{t=1}^{T} \mathcal{F}_{t}\right) \geq 1 - 2\delta.$$

Proof of Lemma 1. Denote  $L_{2,t}(f) = \sum_{i=1}^{M} \sum_{s=1}^{t} |f^{(i)}(\boldsymbol{x}_{s,i}) - y_{s,i}|^2$  and  $\tilde{f}_t = \hat{f}_t - f_\theta$ , we have

$$L_{2,t}(\hat{f}) - L_{2,t}(f_{\theta}) = \sum_{i=1}^{M} \sum_{s=1}^{t} \left| \hat{f}_{t}^{(i)}(\boldsymbol{x}_{s,i}) - y_{s,i} \right|^{2} - \left| f_{\theta}^{(i)}(\boldsymbol{x}_{s,i}) - y_{s,i} \right|^{2}$$
(3)

$$= \sum_{i=1}^{M} \sum_{s=1}^{t} \left| \hat{f}_{t}^{(i)}(\boldsymbol{x}_{s,i}) - f_{\theta}^{(i)}(\boldsymbol{x}_{s,i}) - \eta_{s,i} \right|^{2} - \eta_{s,i}^{2}$$
(4)

$$= \left\| \hat{f}_t - f_\theta \right\|_{2, E_t}^2 - \sum_{i=1}^M \sum_{s=1}^t 2\eta_{s,i} \cdot \tilde{f}_t^{(i)}(\boldsymbol{x}_{s,i})$$
 (5)

By the optimality of  $\hat{f}$ , we know (5)  $\leq$  0, hence

$$\left\| \hat{f}_t - f_\theta \right\|_{2, E_t}^2 \le \sum_{i=1}^M 2 \left\langle \boldsymbol{\eta}_{t, i}, \tilde{f}_t^{(i)}(\boldsymbol{X}_{t, i}) \right\rangle \tag{6}$$

here  $\tilde{f}_t^{(i)}(\boldsymbol{X}_{t,i}) = [\tilde{f}_t^{(i)}(\boldsymbol{x}_{1,i}), \tilde{f}_t^{(i)}(\boldsymbol{x}_{2,i}), \dots, \tilde{f}_t^{(i)}(\boldsymbol{x}_{t,i})]^{\top}$  and  $\boldsymbol{\eta}_{t,i} = [\eta_{1,i}, \eta_{2,i}, \dots, \eta_{t,i}]^{\top}$  are both in  $\mathbb{R}^t$ . We can represent each function  $\tilde{f}_t^{(i)}(\cdot)$  in form  $\tilde{f}_t^{(i)}(\cdot) = \left[\phi^{\star}(\cdot)^{\top}, \hat{\phi}_t(\cdot)^{\top}\right] \begin{bmatrix} \boldsymbol{w}_{t,i}^{\star} \\ -\hat{\boldsymbol{w}}_{t,i} \end{bmatrix} = \phi^{\star}(\cdot)^{\top} \boldsymbol{w}_{t,i}^{\star} - \hat{\phi}_t(\cdot)^{\top} \hat{\boldsymbol{w}}_{t,i}$ , which is exactly  $f_{\theta} - \hat{f}_t$ . Denote  $\tilde{\phi}_t(\cdot) = \begin{bmatrix} \phi^{\star}(\cdot) \\ \hat{\phi}_t(\cdot) \end{bmatrix} \in \Phi^2$  and  $\tilde{\boldsymbol{w}}_{t,i} = \begin{bmatrix} \boldsymbol{w}_{t,i}^{\star} \\ -\hat{\boldsymbol{w}}_{t,i} \end{bmatrix} \in \mathbb{R}^{2k}$ , then  $\tilde{f}_t^{(i)}(\cdot) = \tilde{\phi}_t(\cdot)^{\top} \tilde{\boldsymbol{w}}_{t,i}$ . Since the output of  $\tilde{\phi}_t(\boldsymbol{x}_{s,i}) \in \mathbb{R}^{2k}$ , we can take following decomposition for each  $i \in [M]$ 

$$\tilde{\phi}_t(\boldsymbol{X}_{t,i}) = \left[\tilde{\phi}_t(\boldsymbol{x}_{s,i})\right]_{s=1}^t, \quad \tilde{\phi}_t(\boldsymbol{X}_{t,i})^\top = \boldsymbol{U}_i\boldsymbol{Q}_i, \quad \boldsymbol{U}_i \in \mathcal{O}^{t \times 2k}, \boldsymbol{Q}_i \in \mathbb{R}^{2k \times 2k}.$$

For regret bound, we only need to care about  $t \ge 2k$  by a constant regret difference, hence this decomposition is possible. Plug it into (6) and we get

$$\frac{1}{2} \left\| \hat{f} - f_{\theta} \right\|_{2, E_{t}}^{2} \leq \sum_{i=1}^{M} \left\langle \boldsymbol{\eta}_{t, i}, \tilde{f}_{t}^{(i)}(\boldsymbol{X}_{t, i}) \right\rangle \tag{7}$$

$$= \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot \tilde{\phi}_{t}(\boldsymbol{X}_{t,i})^{\top} \tilde{\boldsymbol{w}}_{t,i}$$
 (8)

$$= \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot \boldsymbol{U}_{i} \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i}$$
 (9)

Notice that, however,  $U_t$  is obtained from optimization problem, which further depends on concrete sampled noise  $\eta_{t,i}$ , hence the concentration bound based on i.i.d. assumption cannot be applied directly. If we fix function  $\tilde{f}_t = \bar{f}_t$ , which induces corresponding  $\bar{\phi}_t(\cdot)$  and  $\bar{\phi}_t(\boldsymbol{X}_{t,i}) = \bar{\boldsymbol{U}}_i(\bar{\phi})\bar{\boldsymbol{Q}}_i$ ,  $\bar{\boldsymbol{U}}_i(\bar{\phi})$  means  $\bar{\boldsymbol{U}}_i$  is a function determined by  $\bar{\phi}$ . According to standard sub-exponential random variable concentration bound, each  $\bar{\boldsymbol{U}}_i(\bar{\phi})$  has 2k independent degrees of freedom, hence we know that with probability at least  $1-\delta_1$ 

$$\sum_{i=1}^{M} \|\bar{\boldsymbol{U}}_{i}^{\top} \boldsymbol{\eta}_{t,i}\|^{2} \le 2Mk + \log(1/\delta_{1})$$
(10)

Denote  $\Phi^2 = \{g(\boldsymbol{x}) = [\phi_1(\boldsymbol{x})^\top, \phi_2(\boldsymbol{x})^\top]^\top : \phi_1, \phi_2 \in \Phi\}, \Phi_\alpha^2$  is an  $\alpha$ -cover of  $\Phi^2$  such that for any  $\phi \in \Phi^2$ , there is a  $\phi_\alpha \in \Phi_\alpha^2$  such that

$$\max_{\boldsymbol{x} \in \mathcal{C} \times \mathcal{A}} \|\phi(\boldsymbol{x}) - \phi_{\alpha}(\boldsymbol{x})\|_{2} \le \alpha.$$
 (11)

For  $\tilde{\phi}$ , find a closest  $\bar{\phi} \in \Phi_{\alpha}^2$  from  $\alpha$ -cover net to satisfy the requirement above, then denote  $\bar{f}_t^{(i)}(\cdot) = \bar{\phi}(\cdot)^{\top} \tilde{\boldsymbol{w}}_{t,i}$ . By union bound, we know that with probability at least  $1 - |\Phi_{\alpha}^2| \delta_1$ , for any  $\bar{\phi} \in \Phi_{\alpha}^2$ , the induced  $\bar{\boldsymbol{U}}_i(\bar{\phi})$  satisfy inequality (10), therefore

$$\frac{1}{2} \left\| \hat{f}_t - f_\theta \right\|_{2, E_t}^2 \le \sum_{i=1}^M \left\langle \boldsymbol{\eta}_{t, i}, \tilde{f}_t^{(i)}(\boldsymbol{X}_{t, i}) \right\rangle \tag{12}$$

$$= \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot \boldsymbol{U}_{i} \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i} = \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot (\boldsymbol{U}_{i} - \bar{\boldsymbol{U}}_{i} + \bar{\boldsymbol{U}}_{i}) \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i}$$
(13)

$$= \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot \bar{\boldsymbol{U}}_{i} \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i} + \sum_{i=1}^{M} \boldsymbol{\eta}_{t,i}^{\top} \cdot (\boldsymbol{U}_{i} - \bar{\boldsymbol{U}}_{i}) \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i}$$
(14)

$$\leq \sqrt{\sum_{i=1}^{M} \left\| \bar{\boldsymbol{U}}_{i}^{\top} \boldsymbol{\eta}_{t,i} \right\|^{2}} \cdot \sqrt{\sum_{i=1}^{M} \left\| \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i} \right\|^{2}} + \sum_{i=1}^{M} \left\langle \boldsymbol{\eta}_{t,i}, \tilde{f}_{t} - \bar{f}_{t} \right\rangle$$
(15)

$$\leq \sqrt{\sum_{i=1}^{M} \left\| \bar{\boldsymbol{U}}_{i}^{\top} \boldsymbol{\eta}_{t,i} \right\|^{2}} \cdot \sqrt{\sum_{i=1}^{M} \left\| \boldsymbol{U}_{i} \boldsymbol{Q}_{i} \tilde{\boldsymbol{w}}_{t,i} \right\|^{2}} + \sum_{i=1}^{M} \left\langle \boldsymbol{\eta}_{t,i}, \tilde{f}_{t} - \bar{f}_{t} \right\rangle$$
(16)

$$= \sqrt{\sum_{i=1}^{M} \left\| \bar{\boldsymbol{U}}_{i}^{\top} \boldsymbol{\eta}_{t,i} \right\|^{2} \cdot \left\| \tilde{\boldsymbol{f}} \right\|_{2,E_{t}} + \sum_{i=1}^{M} \left\langle \boldsymbol{\eta}_{t,i}, \tilde{f}_{t} - \bar{f}_{t} \right\rangle}$$

$$\tag{17}$$

$$\leq \sqrt{2Mk + \log(1/\delta_1)} \cdot \left\| \tilde{f} \right\|_{2, E_t} + \sqrt{\sum_{i=1}^M \|\boldsymbol{\eta}_{t, i}\|^2} \cdot \left\| \tilde{f}_t - \bar{f}_t \right\|_{2, E_t}$$
 (18)

The first term of (18) comes from (10), and the second term is from Cauchy inequality. We assign  $\delta_t = \frac{\delta_2}{T}$  failure probability for event

$$\omega_t : \sum_{i=1}^{M} \|\boldsymbol{\eta}_{t,i}\|^2 \ge Mt + \log(2Mt/\delta_t).$$

By union bound, we have

$$\mathbb{P}\left(\exists t \in [T] : \sum_{i=1}^{M} \|\boldsymbol{\eta}_{t,i}\|^2 \ge Mt + \log(2Mt^2/\delta_2)\right) \le \sum_{t=1}^{T} \delta_t \le \delta_2.$$
 (19)

Next we will give a bound for  $\|\tilde{f}_t - \bar{f}_t\|_{2,E_t}$ 

$$\left\| \tilde{f}_t - \bar{f}_t \right\|_{2, E_t}^2 = \sum_{i=1}^M \sum_{s=1}^t \left| \tilde{\phi}_t(\boldsymbol{x}_{s,i})^\top \tilde{\boldsymbol{w}}_{s,i} - \bar{\phi}_t(\boldsymbol{x}_{s,i})^\top \tilde{\boldsymbol{w}}_{s,i} \right|^2$$
(20)

$$= \sum_{i=1}^{M} \sum_{s=1}^{t} \left| (\tilde{\phi}_t(\boldsymbol{x}_{s,i}) - \bar{\phi}_t(\boldsymbol{x}_{s,i}))^{\top} \tilde{\boldsymbol{w}}_{s,i} \right|^2$$
 (21)

$$\leq \sum_{i=1}^{M} \sum_{s=1}^{t} \left\| \tilde{\phi}_{t}(\boldsymbol{x}_{s,i}) - \bar{\phi}_{t}(\boldsymbol{x}_{s,i}) \right\|_{2}^{2} \cdot \left\| \tilde{\boldsymbol{w}}_{s,i} \right\|_{2}^{2}$$
(22)

According to our assumption, we know  $\|\tilde{\boldsymbol{w}}_{s,i}\|^2 \leq 2\|\boldsymbol{w}_{s,i}\|^2 + 2\|\hat{\boldsymbol{w}}_{s,i}\|^2 \leq 4k$ , from (11) we know  $\|\tilde{\phi}_t(\boldsymbol{x}_{s,i}) - \bar{\phi}_t(\boldsymbol{x}_{s,i})\|_2 \leq \alpha$ , hence

$$\left\| \tilde{f}_t - \bar{f}_t \right\|_{2, E_t}^2 \le 4Mtk\alpha^2 \tag{23}$$

Plug (19) and (23) back into (18), we know with probability at least  $1 - \delta_2 - |\Phi_{\alpha}^2| \delta_1$ , for any  $t \in \mathbb{N}$ 

$$\frac{1}{2} \left\| \tilde{f}_t \right\|_{2,E_t}^2 \le \sqrt{2Mk + \log(1/\delta_1)} \cdot \left\| \tilde{f}_t \right\|_{2,E_t} + \sqrt{Mt + \log(2Mt^2/\delta_2)} \cdot \sqrt{4Mtk\alpha^2} \tag{24}$$

Some simple algebraic transform gives

$$\left\| \hat{f}_t - f_\theta \right\|_{2, E_t}^2 = \left\| \tilde{f}_t \right\|_{2, E_t}^2 \le 6(2Mk + \log(1/\delta_1)) + 8\alpha \sqrt{Mtk(Mt + \log(2Mt^2/\delta_2))}$$
 (25)

Let  $\delta_1 = \delta/|\Phi_{\alpha}^2|, \delta_2 = \delta$ , and notice  $\log |\Phi_{\alpha}^2| \leq 2\log (\mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty}))$ , we conclude that with probability at least  $1-2\delta$ , for every  $t \in \mathbb{N}$ 

$$\left\|\hat{f}_t - f_\theta\right\|_{2E_*}^2 \le 12Mk + 12\log\left(\mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty})/\delta\right) + 8\alpha\sqrt{Mtk(Mt + \log(2Mt^2/\delta))} \quad (26)$$

where the right handside is exactly our defined  $\beta_t(\Phi, \alpha, \delta)$ , hence our conclusion holds.

**Lemma 2.** If  $(\beta_t \geq 0 \mid t \in \mathbb{N})$  is a nondecreasing sequence and  $\mathcal{F}_t := \{f \in \mathcal{F}^{\otimes M} : \|f - \hat{f}_t^{LS}\|_{2,E_t} \leq \sqrt{\beta_t} \}$ . Also, denote  $\mathcal{F} = \mathcal{L} \circ \Phi : \mathcal{C} \times \mathcal{A} \mapsto [0,1]$ , we have

$$\sum_{t=1}^{T} \mathbb{I}\left(w_{\mathcal{F}_t}(\boldsymbol{X}_t) > \epsilon\right) \le \left(\frac{4M\beta_T}{\epsilon^2} + 1\right) \dim_E(\mathcal{F}, \epsilon)$$

*Proof.* The main structure of this proof is similar to proposition 3, section C in Eluder dimension's paper, and we will only point out the subtle details that makes the difference. We will show that if  $w_{\mathcal{F}_t}(\boldsymbol{X}_t) > \epsilon$ , then  $\boldsymbol{X}_t$  is  $\epsilon$ -dependent on fewer than  $4M\beta_T/\epsilon^2$  disjoint subsequences of  $(\boldsymbol{X}_1,\ldots,\boldsymbol{X}_{t-1})$ . Note that if  $w_{\mathcal{F}_t}(\boldsymbol{X}_t) > \epsilon$ , there are  $\overline{f},\underline{f} \in \mathcal{F}_t$  such that  $\sum_{i=1}^M \overline{f}^{(i)}(\boldsymbol{x}_{t,i}) - \underline{f}^{(i)}(\boldsymbol{x}_{t,i}) > \epsilon$ . By definition, if  $\boldsymbol{X}_t$  is  $\epsilon$ -dependent on a subsequence  $(\boldsymbol{X}_{t_1},\boldsymbol{X}_{t_2},\ldots,\boldsymbol{X}_{t_k})$  of  $(\overline{\boldsymbol{X}}_1,\ldots,\boldsymbol{X}_{t-1})$ , then we know

$$\sum_{i=1}^k \left(\sum_{i=1}^M \overline{f}^{(i)}(\boldsymbol{x}_{t_j,i}) - \underline{f}^{(i)}(\boldsymbol{x}_{t_j,i})\right)^2 > \epsilon^2$$

It follows that, if  $X_t$  is  $\epsilon$ -dependent on K disjoint subsequences of  $(X_1, \dots, X_{t-1})$ , then

$$\|\overline{f} - \underline{f}\|_{2,E_{t}}^{2} = \sum_{s=1}^{t} \sum_{i=1}^{M} \left(\overline{f}^{(i)}(\boldsymbol{x}_{s,i}) - \underline{f}^{(i)}(\boldsymbol{x}_{s,i})\right)^{2}$$

$$\geq \frac{1}{M} \sum_{s=1}^{t} \left(\sum_{i=1}^{M} \overline{f}^{(i)}(\boldsymbol{x}_{s,i}) - \underline{f}^{(i)}(\boldsymbol{x}_{s,i})\right)^{2}$$
(Cauchy Inequality)
$$> \frac{K\epsilon^{2}}{M}$$
(28)

By triangle inequality we have

$$\|\overline{f} - \underline{f}\|_{2, E_t} \le \|\overline{f} - \hat{f}_t^{LS}\|_{2, E_t} + \|\hat{f}_t^{LS} - \underline{f}\|_{2, E_t} \le 2\sqrt{\beta_t} \le 2\sqrt{\beta_T}$$
 (29)

and it follows that  $K < 4M\beta_T/\epsilon^2$ .

Notice that essentially we are analyzing scalar output function  $g(\boldsymbol{X}_t) = \sum_{i=1}^M f^{(i)}(\boldsymbol{x}_{t,i})$  where  $f \in \mathcal{F}^{\otimes M}$ . Hence if we denote any  $f \in \mathcal{F}^{\otimes M}$  as  $f(\cdot) = \phi(\cdot)^{\top} \boldsymbol{\Theta}$ , then  $g(\cdot) = \phi(\cdot)^{\top} \boldsymbol{w} \in \mathcal{F}, \boldsymbol{w} = \boldsymbol{\Theta} \cdot \boldsymbol{1}$ . Hence from original eluder dimension paper we know in any action sequence  $(\boldsymbol{X}_1, \dots, \boldsymbol{X}_{\tau})$ , there must exist some element  $\boldsymbol{X}_j$  that is  $\epsilon$ -dependent on at least  $\tau/d-1$  disjoint subsequences of  $(\boldsymbol{X}_1, \dots, \boldsymbol{X}_{\tau})$ , where  $d := \dim_E(\mathcal{F}, \epsilon)$ . Finally we select  $\boldsymbol{X}_1, \dots, \boldsymbol{X}_{\tau}$  as those actions that  $w_{\mathcal{F}_t} > \epsilon$ , combine these two facts above and get  $\tau/d-1 \leq 4M\beta_T/\epsilon^2$ . Hence  $\tau \leq (4M\beta_T/\epsilon^2+1)d$ , which is our desired conclusion.

## **B** Linear MDP Regret Analysis

Apart from the notations section 3, we add more symbols for the regret analysis. We use Q[f] or  $Q[\phi \circ \theta]$  to denote the Q-value function parametrized by function f as Q[f](s,a) = f(s,a) or  $Q[\phi \circ \theta](s,a) = \phi(s,a)^{\top} \theta$  (similar for V[f] as state's value estimation function). Also, based on assumption 2.1, for any  $\left\{Q_{h+1}^{(i)}\right\}_{i=1}^{M}$ , there always exists  $\dot{f}_h\left[Q_{h+1}\right] \in \mathcal{F}^{\otimes M}$  such that

$$\Delta_h^{(i)} \left( Q_{h+1}^{(i)} \right) (s, a) = \mathcal{T}_h^i \left( Q_{h+1}^{(i)} \right) (s, a) - \dot{f}_h^{(i)} (s, a) \tag{30}$$

where the approximation error  $\left\|\Delta_h^{(i)}\left(Q_{h+1}^{(i)}\right)\right\| \leq \mathcal{I}$  for  $\forall i \in [M]$ . Here  $\dot{f}_h[Q_{h+1}]$  indicates that function  $\dot{f}_h$  has dependence on Q-value function  $Q_{h+1}$  on next level h+1. In following analysis, we will use different annotations for different function approximation as below

- $f_h^{(i)*}(\cdot,\cdot) = \phi^*(\cdot,\cdot)^\top \boldsymbol{\theta}_h^{(i)*}$  is the "best" Q-value function approximation in  $\mathcal{Q}_h$  for task i at level h.
- $\hat{f}_h^{(i)}(\cdot,\cdot) = \hat{\phi}(\cdot,\cdot)^{\top}\hat{\theta}_i$  is the empirical least-square minimizer solution for task i at level h.
- $\dot{f}_h^{(i)}(\cdot,\cdot) = \dot{\phi}(\cdot,\cdot)^{\top}\dot{\theta}_i$  is the value approximation function  $\mathcal{T}_h^{(i)}Q_{h+1}^{(i)}$  induced by  $Q_{h+1}^{(i)}$  for task i at level h.
- $\tilde{f}_h^{(i)}(\cdot,\cdot) = \tilde{\phi}(\cdot,\cdot)^{\top}\tilde{\boldsymbol{\theta}}_i$  is the optimism Q-value approximation function for task i at level h.
- $\bar{f}_h^{(i)}(\cdot,\cdot) = \bar{\phi}(\cdot,\cdot)^{\top}\bar{\theta}_i$  is the nearest neighbor in covering set for task i at level h.

#### **B.1** Main Proof sketch

The overall structure is similar to bandits, the main difference here is that we need to take care of the transition dynamics.

Firstly, we decompose the total regret into following terms

$$\operatorname{Reg}(T) = \sum_{t=1}^{T} \sum_{i=1}^{M} \left( V_1^{(i)*} - V_1^{\pi_t^i} \right) \left( s_{1,t}^{(i)} \right)$$
 (31)

$$= \sum_{t=1}^{T} \sum_{i=1}^{M} \left( V_{1}^{(i)\star} - V_{1}^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] \right) \left( s_{1,t}^{(i)} \right) + \sum_{t=1}^{T} \sum_{i=1}^{M} \left( V_{1}^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] - V_{1}^{\pi_{t}^{i}} \right) \left( s_{1,t}^{(i)} \right)$$
(32)

$$\leq \sum_{t=1}^{T} \sum_{i=1}^{M} \left( V_{1}^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] - V_{1}^{\pi_{t}^{i}} \right) \left( s_{1,t}^{(i)} \right) + MHT\mathcal{I}. \tag{33}$$

The inequality is because according to lemma 3, we have at each episode  $t \in [T]$ 

$$\sum_{i=1}^{M} \left( V_1^{i\star} - V_1^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] \right) \left( s_{1,t}^{(i)} \right) \leq MH\mathcal{I}$$

$$\implies \sum_{t=1}^{T} \sum_{i=1}^{M} \left( V_1^{i\star} - V_1^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] \right) \left( s_{1,t}^{(i)} \right) \leq MHT\mathcal{I}.$$

Denote  $a_{h,t}^{(i)} = \pi_t^i \left( s_{ht}^{(i)} \right)$ ,  $Q_h^{(i)}[\tilde{f}_{h,t}^{(i)}] = \tilde{Q}_{h,t}^{(i)}$  and  $V_h^{(i)}[\tilde{f}_{h,t}^{(i)}] = \tilde{V}_{h,t}^{(i)}$  for short. We have for any  $t \in [T], h \in [H]$ 

$$\sum_{i=1}^{M} \left( \tilde{V}_{h,t}^{(i)} - V_{h,t}^{\pi_{t}^{i}} \right) \left( s_{h,t}^{(i)} \right) = \sum_{i=1}^{M} \left( \tilde{Q}_{h,t}^{(i)} - Q_{h,t}^{\pi_{t}^{i}} \right) \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \\
= \sum_{i=1}^{M} \left( \tilde{Q}_{h,t}^{(i)} - \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} \right) \left( s_{1,t}^{(i)}, a_{h,t}^{(i)} \right) + \sum_{i=1}^{M} \left( \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} - Q_{h,t}^{\pi_{t}^{i}} \right) \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \\
(35)$$

Since the failure event  $\bigcup_{t=1}^T \bigcup_{h=1}^H E_{ht}$  only happens with probability  $\delta$  according to lemma 6, and the addition of regret when it happens is constant bounded, we will simply assume that it does not happen. Then applying lemma 5, we have

$$\sum_{i=1}^{M} \left( \tilde{Q}_{h,t}^{(i)} - \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} \right) \left( s_{h,t}^{(i)}, \ a_{h,t}^{(i)} \right) \le M\mathcal{I} + 2w_{\mathcal{F}_{h,t}} \left( \boldsymbol{x}_{h,t} \right). \tag{36}$$

where  $\boldsymbol{x}_{h,t} = \left[ (s_{h,t}^{(1)}, a_{h,t}^{(1)}), \dots, (s_{h,t}^{(M)}, a_{h,t}^{(M)}) \right]$  denotes the stacked input for all state-action pair at level h, episode t.

Next, we expand the second summation in (35) and have

$$\sum_{i=1}^{M} \left( \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} - Q_{h,t}^{\pi_{t}^{i}} \right) \left( s_{h,t}^{(i)}, \ a_{h,t}^{(i)} \right) = \sum_{i=1}^{M} \mathbb{E}_{s' \sim \mathcal{P}_{h}^{(i)} \left( \cdot \mid s_{h,t}^{(i)}, a_{h,t}^{(i)} \right)} \left[ \left( \tilde{V}_{h+1,t}^{(i)} - V_{h+1}^{\pi_{t}^{i}} \right) \left( s' \right) \right]$$

$$= \sum_{i=1}^{M} \left( \tilde{V}_{h+1,t}^{(i)} - V_{h+1}^{\pi_{t}^{i}} \right) \left( s_{h+1,t}^{(i)} \right) + \sum_{i=1}^{M} \zeta_{h,t}^{(i)}$$
(38)

where  $\zeta_{h,t}^{(i)}$  is a martingale difference with respect to history  $\mathcal{H}_{h,t}$  defined by

$$\zeta_{h,t}^{(i)} \stackrel{\text{def}}{=} \mathbb{E}_{s' \sim \mathcal{P}_h^{(i)}\left(\cdot \mid s_{h,t}^{(i)}, a_{h,t}^{(i)}\right)} \left[ \left( \tilde{V}_{h+1,t}^{(i)} - V_{h+1}^{\pi_t^i} \right) \left( s' \right) \right] - \left( \tilde{V}_{h+1,t}^{(i)} - V_{h+1}^{\pi_t^i} \right) \left( s' \right)$$
(39)

According to assumption 2.2 we know that  $|\zeta_{h,t}^{(i)}| \le 4$ , hence by Azuma-Hoeffding's inequality, we know that with probability at least  $1 - \delta/2$ , for any  $t \in [T]$  and  $i \in [M]$ 

$$\sum_{i=1}^{t} \zeta_{h,t}^{(i)} \le 4\sqrt{2t\log\frac{2T}{\delta}}.\tag{40}$$

We can then apply (38) recursively from h = 1 to H, which gives

$$\operatorname{Reg}(T) \le \sum_{t=1}^{T} \sum_{i=1}^{M} \left( \tilde{V}_{1,t}^{(i)} - V_{1}^{\pi_{t}^{i}} \right) \left( s_{1,t}^{(i)} \right) + MHT\mathcal{I}$$
(41)

$$\leq 2MHT\mathcal{I} + \sum_{t=1}^{T} \sum_{h=1}^{H} 2w_{\mathcal{F}_t}(\boldsymbol{x}_{h,t}) + \sum_{i=1}^{M} \sum_{h=1}^{H} \sum_{t=1}^{T} \zeta_{h,t}^{(i)}$$
(42)

According to lemma 2 we know that

$$\sum_{t=1}^{T} w_{\mathcal{F}_t}(\boldsymbol{x}_{h,t}) \le \left(\frac{4M\beta_{h,T}}{\alpha^2} + 1\right) \dim_E(\mathcal{F}, \alpha)$$
(43)

where  $\beta_{h,t} = \tilde{O}(Mk + \log \mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty}) + MT\mathcal{I}^2)$ . Summarizing all inequality above and we have the final regret bound as

$$\operatorname{Reg}(T) = 2MHT\mathcal{I} + \sum_{t=1}^{T} \sum_{h=1}^{H} 2w_{\mathcal{F}_{t}}(\boldsymbol{x}_{h,t}) + \sum_{i=1}^{M} \sum_{h=1}^{H} \sum_{t=1}^{T} \zeta_{h,t}^{(i)}$$

$$= \tilde{O}\left(MHT\mathcal{I} + \tilde{O}(\sqrt{Mk + \log \mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty}) + MT\mathcal{I}^{2}})H\sqrt{MT \operatorname{dim}_{E}(\mathcal{F}, \alpha)} + MH\sqrt{T}\right)$$
(45)

Set  $\alpha = \frac{1}{kMT}$ , we have the regret bound as

$$\tilde{O}\left(H\sqrt{\dim_E(\mathcal{F},(kMT)^{-1})}\left(M\sqrt{Tk}+\sqrt{MT\log\mathcal{N}(\Phi,(kMT)^{-1},\|\cdot\|_{\infty})}+MT\mathcal{I}\right)\right).$$

## **B.2** Detailed Lemma Proof

**Lemma 3.** Let  $V_1^{i\star}$  be the value of optimal policy and  $V_1^i\left[\tilde{f}_{1,t}^{(i)}\right]$  be the optimistic value estimation defined in main proof. We have the accuracy guarantee as

$$\sum_{i=1}^{M} \left( V_1^{(i)\star} - V_1^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] \right) \left( s_{1,t}^{(i)} \right) \le MH\mathcal{I}. \tag{46}$$

*Proof.* Recursively define the closest value approximator function  $f_h^* = (\phi_h^*)^\top \Theta_h^*$  at level h within function class  $\mathcal{F}^{\otimes M}$  as

$$\phi_h^*, \boldsymbol{\Theta}_h^* \stackrel{\text{def}}{=} \underset{\phi \in \Phi, \boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M] \in \mathbb{R}^{k \times M}}{\arg \min} \underset{s, a, i}{\sup} \left| \phi(s, a)^\top \boldsymbol{\theta}_h^{(i)} - \mathcal{T}_h^{(i)} Q_{h+1}^{(i)} \left[ \phi_{h+1}^* \circ \boldsymbol{\theta}_{h+1}^{(i)*} \right] (s, a) \right| \quad (47)$$

with  $\boldsymbol{\theta}_{H+1}^{(i)} = \mathbf{0}$  for any  $i \in [M]$  and  $\boldsymbol{\Theta}_h^* = \left[\boldsymbol{\theta}_h^{(1)*}, \dots, \boldsymbol{\theta}_h^{(M)*}\right]$ . By lemma 6 in [44] we have

$$\sup_{(s,a)\in\mathcal{S}\times\mathcal{A},i\in[M]} \left| Q_h^{(i)\star}(s,a) - \phi_h^*(s,a)^\top \boldsymbol{\theta}_h^{(i)\star} \right| \le (H-h+1)\mathcal{I}. \tag{48}$$

where  $Q_h^{(i)\star}$  is the optimal value function for task i.

Next, we will show that  $f_h^*$  is a feasible solution for the optimization of  $\mathcal{F}_t$ . This is achieved via inductive construction. For h=H+1 we know it holds trivially because  $\tilde{f}_{H+1}^{(i)}=f_{H+1}^{(i)*}=\mathbf{0}$ . Now we suppose that  $\beta_{h,t}$  for  $k=h+1,\ldots,H$  satisfies that we can always find  $\tilde{f}_k^{(i)}=f_k^{(i)*}$ . Then from the definition of  $f_h^{(i)*}$  we can always properly set  $\mathcal{F}_{h,t}$  (to be specified later) to let it contain

$$\dot{f}_{h}^{(i)} \left[ V_{h+1}^{(i)} \left[ f_{h+1}^{(i)*} \right] \right] = f_{h}^{(i)*}. \tag{49}$$

By lemma 4, we have

$$\left\| \hat{f}_h \left[ V_{h+1} \left[ f_{h+1}^* \right] \right] - \dot{f}_h \left[ V_{h+1} \left[ f_{h+1}^* \right] \right] \right\|_{2, E_t}^2 \le \beta_{h, t}. \tag{50}$$

Therefore, set  $\beta_{h,t}$  as the function we set *does* let  $f_h^{(i)*} \in \mathcal{F}_{h,t}$ .

Finally, we can finish the proof from showing that

$$\sum_{i=1}^{M} V_1^{(i)} \left[ \tilde{f}_{1,t}^{(i)} \right] \left( s_{1,t}^{(i)} \right) \tag{51}$$

$$= \sum_{i=1}^{M} \max_{a \in \mathcal{A}} \tilde{f}_{1,t}^{(i)} \left( s_{1,t}^{(i)}, a \right)$$
 (52)

$$\geq \sum_{i=1}^{M} \max_{a \in \mathcal{A}} f_{1,t}^{(i)*} \left( s_{1,t}^{(i)}, a \right) \qquad \qquad \text{(because } f_1^{(i)*} \in \mathcal{F}_t \text{)}$$

$$\geq \sum_{i=1}^{M} f_{1,t}^{(i)*} \left( s_{1,t}^{(i)}, \pi_{1}^{i*} \left( s_{1,t}^{(i)} \right) \right) \tag{53}$$

$$\geq \sum_{i=1}^{M} Q_{1}^{(i)\star} \left( s_{1,t}^{(i)}, \pi_{1}^{i\star} \left( s_{1,t}^{(i)} \right) \right) - MH\mathcal{I}$$
 (By (48))

$$\geq \sum_{i=1}^{M} V_1^{(i)\star} \left( s_{1,t}^{(i)} \right) - MH\mathcal{I}. \tag{54}$$

**Lemma 4.** For any episode  $t \in [T]$ , level  $h \in [H]$  and any Q-value function at next level  $\{Q_{h+1}^{(i)}\}_{i=1}^M \in \mathcal{Q}_{h+1}$ , denote  $\dot{f}_{h,t}$  as the best fit Q-value estimation induced by  $Q_{h+1}^{(i)}$  minimizing Bellman error, we have

$$\left\| \hat{f}_{h,t} \left[ Q_{h+1} \right] - \dot{f}_{h,t} \left[ Q_{h+1} \right] \right\|_{2,E_t}^2 \le \beta_{h,t} \stackrel{\text{def}}{=} \left( B_{h,1} + \sqrt{MT} \mathcal{I} + \sqrt{B_{h,2}} \right)^2. \tag{55}$$

The  $B_{h,1}$  and  $B_{h,2}$  are from Lemma 6. Equivalently saying, this means that  $\dot{f}_{h,t}$  is contained in set  $\mathcal{F}_{h,t}$  defined as

$$\mathcal{F}_{h,t} \stackrel{\text{def}}{=} \left\{ f \in \mathcal{F}^{\otimes M} : \left\| f - \hat{f}_{h,t} \left[ Q_{h+1} \right] \right\|_{2,E_t}^2 \le \beta_{h,t} \right\}.$$

*Proof.* By the empirical optimality of  $\hat{f}_{h,t}$ , we know

$$\sum_{i=1}^{M} \left\| \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \boldsymbol{y}_{h,t}^{(i)} \right\|^{2} \le \sum_{i=1}^{M} \left\| \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \boldsymbol{y}_{h,t}^{(i)} \right\|^{2}.$$
 (56)

Here we abuse the notation and use  $\hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t})$  to denote function  $\hat{f}_{h,t}^{(i)}$ 's output on all the state-action pair  $\boldsymbol{X}_{h,t}$  in the first t-1 episodes at level h for task i, also  $\boldsymbol{y}_{h,t}^{(i)}$  is the corresponding target value label. This inequality implies that

$$\sum_{i=1}^{M} \left\| \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\|^{2}$$
(57)

$$\leq 2\sum_{i=1}^{M} \left\langle \boldsymbol{\Delta}_{h,t}^{(i)}, \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle + 2\sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)}, \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
(58)

where

$$\boldsymbol{\Delta}_{h,t}^{(i)} \stackrel{\text{def}}{=} \begin{bmatrix} \Delta_{h,1}^{(i)}(Q_{h+1}^{(i)})(s_{h,1}^{(i)},a_{h,2}^{(i)}) & \Delta_{h,2}^{(i)}(Q_{h+1}^{(i)})(s_{h,2}^{(i)},a_{h,2}^{(i)}) & \dots & \Delta_{h,t-1}^{(i)}(Q_{h+1}^{(i)})(s_{h,t-1}^{(i)},a_{h,t-1}^{(i)}) \end{bmatrix}$$

is the Bellman error for Q-value approximation, each  $\Delta_{h,j}^{(i)}(Q_{h+1}^{(i)})(s_{h,j}^{(i)},a_{h,j}^{(i)})$  is defined in (30). And

$$\boldsymbol{z}_{h,t}^{(i)} \stackrel{\text{def}}{=} \left[ z_{h,1}^{(i)}(Q_{h+1}^{(i)})(s_{h,1}^{(i)}, a_{h,2}^{(i)}) \quad \dots \quad z_{h,t-1}^{(i)}(Q_{h+1}^{(i)})(s_{h,t-1}^{(i)}, a_{h,t-1}^{(i)}) \right]$$

where 
$$z_{h,j}^{(i)}\left(Q_{h+1}^{(i)}\right)\left(s_{h,j}^{(i)},a_{h,j}^{(i)}\right) \stackrel{\text{def}}{=} R\left(s_{h,j}^{(i)},a_{h,j}^{(i)}\right) + \max_{a\in\mathcal{A}}Q_{h+1}^{(i)}\left(s_{h+1,j}^{(i)},a\right) - \mathcal{T}_h^{(i)}\left(Q_{h+1}^{(i)}\right)\left(s_{h,j}^{(i)},a_{h,j}^{(i)}\right) \text{ is the finite sampling noise.}$$

Next, we are going to bound the two terms in (58). For the first term, we have

$$\sum_{i=1}^{M} \left\langle \Delta_{h,t}^{(i)}, \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
 (59)

$$\leq \sum_{i=1}^{M} \left\| \boldsymbol{\Delta}_{h,t}^{(i)} \right\| \cdot \left\| \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\|$$
(60)

$$\leq \sqrt{T}\mathcal{I} \cdot \sum_{i=1}^{M} \left\| \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\|$$
(61)

$$\leq \sqrt{MT} \mathcal{I} \cdot \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2,E_t} \tag{62}$$

By lemma 6, when the failure case does not happen, we have

$$\sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)}, \hat{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h,t}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle \leq B_{h,1} \cdot \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2,E_{t}} + B_{h,2}$$
 (63)

where

$$B_{h,1} = \sqrt{2Mk + \log(\mathcal{N}(\Phi, (kMT)^{-1}, \|\cdot\|_{\infty})/\delta)} + 1$$
(64)

$$B_{h,2} = 2\sqrt{MT + \log(2MT^2/\delta)} \tag{65}$$

Adding the bound for two terms and we get

$$\left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2,E_t}^2 \le (B_{h,1} + \sqrt{MT}\mathcal{I}) \cdot \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2,E_t} + B_{h,2}$$
 (66)

$$\implies \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2}^{2} \le \left( B_{h,1} + \sqrt{MT}\mathcal{I} + \sqrt{B_{h,2}} \right)^{2} \stackrel{\text{def}}{=} \beta_{h,t} \tag{67}$$

which completes the proof.

**Lemma 5.** If the failure event in lemma 6 does not happen, for any feasible solution  $Q_h^{(i)}\left[\tilde{f}_h^{(i)}\right]$  in the definition of  $\mathcal{F}_{h,t}$ , and any  $h\in[H]$ ,  $t\in[T]$ , we have

$$\sum_{i=1}^{M} \left| \left( \tilde{Q}_{h,t}^{(i)} - \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} \right) \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right| \le M \mathcal{I} + 2 w_{\mathcal{F}_{h,t}} \left( \boldsymbol{x}_{h,t} \right), \tag{68}$$

where  $\mathbf{x}_{h,t} = \left[ (s_{h,t}^{(1)}, a_{h,t}^{(1)}), \dots, (s_{h,t}^{(M)}, a_{h,t}^{(M)}) \right]$  denotes the stacked input for all state-action pair at level h, episode t.

Proof.

$$\sum_{i=1}^{M} \left| \left( \tilde{Q}_{h,t}^{(i)} - \mathcal{T}_{h}^{(i)} \tilde{Q}_{h+1,t}^{(i)} \right) \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right| \tag{69}$$

$$= \sum_{i=1}^{M} \left| \tilde{Q}_{h,t}^{(i)}(s,a) - \dot{f}_{h}^{(i)} \left[ \tilde{Q}_{h+1}^{(i)} \right] \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) - \Delta_{h}^{(i)} \left( \tilde{Q}_{h+1}^{(i)} \right) \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right|$$
(70)

$$\leq M\mathcal{I} + \sum_{i=1}^{M} \left| \tilde{f}_{h,t}^{(i)} \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) - \dot{f}_{h}^{(i)} \left[ \tilde{Q}_{h+1}^{(i)} \right] \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right| \tag{71}$$

$$\leq M\mathcal{I} + \sum_{i=1}^{M} \left| \hat{f}_{h,t}^{(i)} \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) - \hat{f}_{h}^{(i)} \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right| + \left| \hat{f}_{h}^{(i)} \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) - \hat{f}_{h}^{(i)} \left[ \tilde{Q}_{h+1}^{(i)} \right] \left( s_{h,t}^{(i)}, a_{h,t}^{(i)} \right) \right|$$

$$(72)$$

According to our construction, we know that both  $\tilde{f}_{h,t}^{(i)}$  and  $\dot{f}_h^{(i)}$  are contained in  $\mathcal{F}_{h,t}$ , therefore we have  $\sum_{i=1}^{M}\left|\tilde{f}_{h,t}^{(i)}\left(s_{h,t}^{(i)},a_{h,t}^{(i)}\right)-\hat{f}_{h}^{(i)}\left(s_{h,t}^{(i)},a_{h,t}^{(i)}\right)\right| \leq w_{\mathcal{F}_{h,t}}\left(x_{h,t}\right)$  and  $\sum_{i=1}^{M}\left|\dot{f}_{h,t}^{(i)}\left[\tilde{Q}_{h+1}^{(i)}\right]\left(s_{h,t}^{(i)},a_{h,t}^{(i)}\right)-\hat{f}_{h}^{(i)}\left(s_{h,t}^{(i)},a_{h,t}^{(i)}\right)\right| \leq w_{\mathcal{F}_{h,t}}\left(x_{h,t}\right)$ , where  $x_{h,t}=\left[\left(s_{h,t}^{(1)},a_{h,t}^{(1)}\right),\ldots,\left(s_{h,t}^{(M)},a_{h,t}^{(M)}\right)\right]$  denotes the stacked input for all state-action pair at level h, episode t.

Summarizing all the inequalities and we know the whole lemma holds.

**Lemma 6.** (Probability bound for failure event) In this lemma we denote  $\hat{f}_h^{(i)}\left[Q_{h+1}^{(i)}\right]$  as  $\hat{f}_h^{(i)}$  for the sake of simplicity (similar for  $\hat{f}_h^{(i)}$ ). Define event  $E_{h,t}$  as

$$E_{h,t} \stackrel{\text{def}}{=} \mathbb{I} \left[ \exists \{ Q_{h+1}^{(i)} \}_{i=1}^{M} \quad \sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle > B_{h,1} \cdot \left\| \hat{f}_{h}^{(i)} - \dot{f}_{h}^{(i)} \right\|_{2,E_{t}} + B_{h,2} \right]$$

$$(73)$$

where  $B_{h,1}$  and  $B_{h,2}$  will be specified later. We have

$$\mathbb{P}\left(\bigcup_{t=1}^{T}\bigcup_{h=1}^{H}E_{h,t}\right) \leq \delta. \tag{74}$$

*Proof.* Similar to lemma 1, we can find a  $\alpha$ -cover  $\Phi_{\alpha}$  for  $\Phi$  such that for any Q-value function  $\left(Q_{h+1}^{(1)}[\phi\circ\theta_1],Q_{h+1}^{(2)}[\phi\circ\theta_2],\dots,Q_{h+1}^{(M)}[\phi\circ\theta_M]\right)$ , we can find  $\bar{\phi}\in\Phi_{\alpha}$  and  $\bar{\theta}_i$  for  $i\in[M]$  such that for any  $(s,a)\in\mathcal{S}\times\mathcal{A}$  and any  $i\in[M]$ 

$$\left| Q_{h+1}^{(i)}(s,a) - \bar{\phi}(s,a)^{\top} \bar{\boldsymbol{\theta}}_i \right| \le \sqrt{k}\alpha. \tag{75}$$

Define  $ar{Q}_{h+1}^{(i)} = Q_{h+1}^{(i)} \left[ ar{\phi} \circ \pmb{\theta}_i 
ight]$  and further let

$$\bar{\boldsymbol{z}}_{h,t}^{(i)} \stackrel{\text{def}}{=} \left[ z_{h,1}^{(i)} \left( \bar{Q}_{h+1}^{(i)} \right) \left( s_{h,1}^{(i)}, a_{h,1}^{(i)} \right) \quad \dots \quad z_{h,t-1}^{(i)} \left( \bar{Q}_{h+1}^{(i)} \right) \left( s_{h,t-1}^{(i)}, a_{h,t-1}^{(i)} \right) \right] \in \mathbb{R}^{t-1}$$

then we have

$$\sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
 (76)

$$= \sum_{i=1}^{M} \left\langle \bar{z}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
 (77)

$$+\sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)} - \bar{\boldsymbol{z}}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
 (78)

(79)

Notice that for fixed  $\bar{f}_h^{(i)}(\cdot,\cdot)=\phi(\cdot,\cdot)^{\top}\bar{\boldsymbol{\theta}}_{h+1}^{(i)}$ , each  $z_{h,1}^{(i)}\left(\bar{Q}_{h+1}^{(i)}\right)\left(s_{h,1}^{(i)},a_{h,2}^{(i)}\right)$  is a zero-mean 1-sub-Gaussian random variable conditioned on past history. Therefore we can treat it as  $\eta_{t,i}=z_{h,t}^{(i)}$  in Lemma 1 and get

$$\sum_{i=1}^{M} \left\langle \bar{z}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
(80)

$$\leq \sqrt{2Mk + \log(1/\delta_1)} \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2.E_t} + 2\alpha \sqrt{Mtk(Mt + \log(2Mt^2/\delta_2))}. \tag{81}$$

Setting  $\delta_1 = \frac{\delta}{2|\Phi^{\alpha}|}, \delta_2 = \delta/2$  and get

$$\sum_{i=1}^{M} \left\langle \bar{z}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
(82)

$$\leq \sqrt{2Mk + \log(\mathcal{N}(\Phi, \alpha, \|\cdot\|_{\infty})/\delta)} \cdot \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2, E_t} + 2\alpha \sqrt{MTk(MT + \log(2MT^2/\delta))}. \tag{83}$$

By union bound, we know it holds for any  $\bar{f}_h$  with probability at least  $1-|\Phi^\alpha|\delta_1=1-\delta$ . Also, from  $\left|Q_{h+1}^{(i)}(s,a)-\bar{\phi}(s,a)^\top\bar{\boldsymbol{\theta}}_i\right|\leq \sqrt{k}\alpha'$  we know that

$$\begin{aligned} & \left| z_{h,j}^{(i)} \left( Q_{h+1}^{(i)} \right) \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right) - z_{h,j}^{(i)} \left( \bar{Q}_{h+1}^{(i)} \right) \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right) \right| \\ & = \left| \max_{a \in \mathcal{A}} Q_{h+1}^{(i)} \left( s_{h+1,j}^{(i)}, a \right) - \mathcal{T}_h^{(i)} \left( Q_{h+1}^{(i)} \right) \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right) - \max_{a \in \mathcal{A}} \bar{Q}_{h+1}^{(i)} \left( s_{h+1,j}^{(i)}, a \right) + \mathcal{T}_h^{(i)} \left( \bar{Q}_{h+1}^{(i)} \right) \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right) \right| \end{aligned} \tag{84}$$

$$\leq \max_{a \in \mathcal{A}} \left| Q_{h+1}^{(i)} \left( s_{h+1,j}^{(i)}, a \right) - \bar{Q}_{h+1}^{(i)} \left( s_{h+1,j}^{(i)}, a \right) \right| + \left| \mathcal{T}_{h}^{(i)} \left( \bar{Q}_{h+1}^{(i)} - Q_{h+1}^{(i)} \right) \left( s_{h,j}^{(i)}, a_{h,j}^{(i)} \right) \right| \tag{86}$$

$$\leq 2\sqrt{k}\alpha'$$
 (87)

hence we have

$$\sum_{i=1}^{M} \left\langle \boldsymbol{z}_{h,t}^{(i)} - \bar{\boldsymbol{z}}_{h,t}^{(i)}, \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\rangle$$
(88)

$$\leq \sum_{i=1}^{M} \left\| \boldsymbol{z}_{h,t}^{(i)} - \bar{\boldsymbol{z}}_{h,t}^{(i)} \right\| \cdot \left\| \hat{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) - \dot{f}_{h}^{(i)}(\boldsymbol{X}_{h,t}) \right\|$$
(89)

$$\leq 2\alpha' \sqrt{MTk} \cdot \left\| \hat{f}_{h,t} - \dot{f}_{h,t} \right\|_{2.E_t} \tag{90}$$

holds for arbitrary  $\{Q_{h+1}^{(i)}\}$  at any level  $h \in [H], t \in [T]$ .

Adding (83) and (90), we finally finish the proof by setting  $\alpha = \alpha' = \frac{1}{MTk}$ 

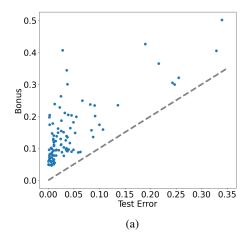
$$B_{h,1} = \sqrt{2Mk + \log(\mathcal{N}(\Phi, (kMT)^{-1}, \|\cdot\|_{\infty})/\delta)} + 1$$
(91)

$$B_{h,2} = 2\sqrt{MT + \log(2MT^2/\delta)} \tag{92}$$

# C Experiment Dissection and Discussion

In this section, we will take a closer view of the learning procedure and analyze the functionality of the UCB term in our algorithm. Usually, a reasonable UCB term should embrace several properties. (i) It should let confidence set  $\mathcal{F}_t$  contain the real parameter with high probability. (ii) It should shrink at a reasonable speed to achieve low regret.

To check (i), we choose the model  $\hat{f}_t$  at step t=200 which is trained on insufficient data with only 2000 samples. We then sample 100 images from test set as unknown inputs  $\mathcal{D}=\{(\boldsymbol{x}_i,y_i)\}_{i=1}^{100}$ , where  $\boldsymbol{x}_i$  is the digit image and  $y_i$  is the corresponding target value. We inspect the relationship between the original prediction error  $|\hat{f}_t(\boldsymbol{x}_i)-y_i|$  and the added bonus  $b_i=\bar{f}_t(\boldsymbol{x}_i)-\hat{f}_t(\boldsymbol{x}_i)$  via finetuning on each input  $\boldsymbol{x}_i\in\mathcal{D}$ . The result is presented as scatter dots in Figure 2(a). We can clearly see that almost all the points lie above the line y=x, meaning that  $b_i=\bar{f}_t(\boldsymbol{x}_i)-\hat{f}_t(\boldsymbol{x}_i)\geq |\hat{f}_t(\boldsymbol{x}_i)-y_i|\geq y_i-\hat{f}_t(\boldsymbol{x}_i)$  for any  $i\in[100]$ , which further indicates that  $\bar{f}_t(\boldsymbol{x}_i)\geq y_i$ . This validates that we can always find some  $\bar{f}\in\mathcal{F}_t$  to give an optimistic estimation of the value for almost every  $\boldsymbol{x}$ . Moreover, we can observe an apparent correlated pattern between the test error and bonus, which implies that our



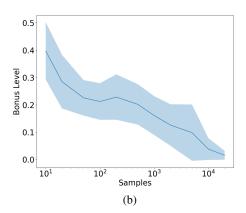


Figure 2: (a) The relationship between unknown data's prediction error and the bonus it gets from finetuning. The grey line is y=x. (b) The average bonus level of 100 test images with respect to the number of samples in training set, the shaded area is the interval for  $\pm 1$  standard deviation.

algorithm will give larger bonus for the data point whose prediction is not reliable, and only give relatively small bonus for the data that it is confident with.

We also check (ii) by plotting the average bonus level (closely related to the width of confidence set) against the number of samples the algorithm has been trained on. We gradually increase the number of samples from 10 to 20000 and fix a set of test images  $\mathcal{D}$  as before to see how the average bonus level changes when the training set size increases. The result is shown in Figure 2(b). Previous work [12] proves that the eluder dimension of neural networks can be exponentially large in the worst case, which means that it can give almost arbitrary output value even when it is constrained to give a precisely accurate prediction for a large number of samples in the training set. In that case, the average bonus level should have remained constant regardless of the size of the training set. However, our experiment shows that the average bonus drops when the number of training samples increases. We conjecture that it is because in reality, when the input data are restricted to regular images with clear semantics, and the optimization procedure of the model is conducted via gradient-based methods in a very close neighborhood, the arbitrariness of the neural network's output is substantially reduced.

Restricting the model's training loss in the training set effectively limits the bonus obtained from the finetune procedure, which realizes the desired fast-shrinking property from our functional confidence set. Such a phenomenon sheds light on the unknown property of neural network's generalization capability and interpolation plasticity. We leave explaining the underlying mechanism as future work.

### C.1 Visualize the Learned Representation

A natural and interesting question is what representation does our CNN backbone actually learn. To investigate this problem and visualize the learned representation, we measure the information of different digits within the learned representation. Interestingly, we find that our model indeed learns an indicative representation for classification problem via multitask value regression training.

The basic measurement for the quality of representation is evaluated with the kernel function  $\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \langle \phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j) \rangle$  and see whether it has a strong diagonal. We take the checkpoint of neural network model at final step (around 600 with more than 6000 samples), and treat the module before the final linear layer as  $\phi(\cdot)$ . Denote the MNIST test set as  $\mathcal{D} = \{\mathcal{D}_i\}_{i=0}^9$  where  $\mathcal{D}_i$  is the images of digit i. Define the correlation between digit i and j under representation  $\phi$  as

$$C(i,j) = \frac{1}{|\mathcal{D}_i| \times |\mathcal{D}_j|} \sum_{\boldsymbol{x}_s \in \mathcal{D}_i} \sum_{\boldsymbol{x}_t \in \mathcal{D}_j} \langle \phi(\boldsymbol{x}_s), \phi(\boldsymbol{x}_t) \rangle$$
(93)

To accelerate the evaluation, notice that we can preprocess an "template vector"  $T_i$  for each digit i as

$$T_i = \frac{1}{|\mathcal{D}_i|} \sum_{\boldsymbol{x} \in \mathcal{D}_i} \phi(\boldsymbol{x})$$
 (94)

so that the correlation can be computed through

$$C(i,j) = \frac{1}{|\mathcal{D}_i| \times |\mathcal{D}_j|} \sum_{\boldsymbol{x}_s \in \mathcal{D}_i} \sum_{\boldsymbol{x}_t \in \mathcal{D}_j} \langle \phi(\boldsymbol{x}_s), \phi(\boldsymbol{x}_t) \rangle$$
(95)

$$= \frac{1}{|\mathcal{D}_j|} \sum_{\boldsymbol{x}_t \in \mathcal{D}_j} \left( \frac{1}{|\mathcal{D}_i|} \sum_{\boldsymbol{x}_s \in \mathcal{D}_i} \langle \phi(\boldsymbol{x}_s), \phi(\boldsymbol{x}_t) \rangle \right)$$
(96)

$$= \frac{1}{|\mathcal{D}_j|} \sum_{\boldsymbol{x}_t \in \mathcal{D}_j} \left\langle \frac{1}{|\mathcal{D}_i|} \sum_{\boldsymbol{x}_s \in \mathcal{D}_i} \phi(\boldsymbol{x}_s), \phi(\boldsymbol{x}_t) \right\rangle$$
(97)

$$= \frac{1}{|\mathcal{D}_j|} \sum_{\boldsymbol{x}_t \in \mathcal{D}_j} \langle \boldsymbol{T}_i, \phi(\boldsymbol{x}_t) \rangle$$
 (98)

$$= \langle \boldsymbol{T}_i, \boldsymbol{T}_i \rangle \tag{99}$$

We plot this 10x10 correlation map for single task training and multitask training with M=10. Notice that the single task reward mapping function is  $\sigma(i)=i/10$ , and to assure the different tasks in multitask training are heterogeneous, we manually set that the best digit for each task are distinct.

The result is in figure 3. We can see that since single task only needs to recognize the large value digit, namely 9, 8 or 7, its representation function is not informative for distinguishing digits. And interestingly, the multitask trained network's representation demonstrates a very strong diagonal, indicating that the representation vector is very specific to the digit's image, although the training process has no explicit definition for the classification task but a regression problem instead. Actually, we found a simple linear layer append to this representation can achieve over 95% accuracy on MNIST test set.

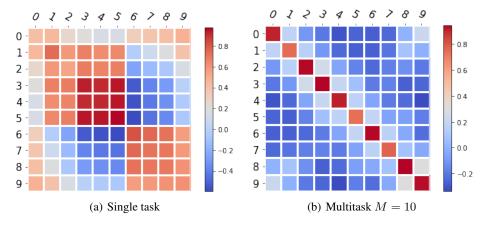


Figure 3: The kernel function for the representation learned by single task and 10-tasks multitask. It is clear that multitask representation learning obtains a more comprehensive and interpretable pattern for the MNIST images.