Iterative Machine Teaching for Black-box Markov Learners

Chaoqi Wang 1 Sandra Zilles 2 Adish Singla 3 Yuxin Chen 1

Abstract

Machine teaching has traditionally been constrained by the assumption of a fixed learner's model. In this paper, we challenge this notion by proposing a novel black-box Markov learner model, drawing inspiration from decision psychology and neuroscience where learners are often viewed as black boxes with adaptable parameters. We model the learner's dynamics as a Markov decision process (MDP) with unknown parameters, encompassing a wide range of learner types studied in machine teaching literature. This approach reduces teaching complexity to finding an optimal policy for the underlying MDP. Building on this, we introduce an algorithm for teaching in this black-box setting and provide an analysis of teaching costs under different scenarios. We further establish a connection between our model and two types of learners in psychology and neuroscience, the epiphany learner and the non-epiphany learner, linking them with discounted and non-discounted black-box Markov learners respectively. This alignment offers a psychologically and neuroscientifically grounded perspective to our work. Supported by numerical study results, this paper delivers a significant contribution to machine teaching, introducing a robust, versatile learner model with a rigorous theoretical foundation.

1. Introduction

Machine teaching seeks effective policies for selecting training examples to help a learner learn a target concept. Over the past few decades, the field of machine teaching has been pushed forward and shown great promise in various application domains, including those targeting human learners, such as automated tutoring systems (Zhu, 2015; Rafferty

Proceedings of the First Workshop on Theory of Mind in Communicating Agents, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

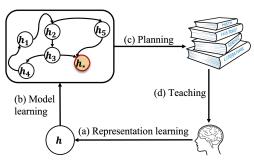


Figure 1: The teaching framework. Our work focuses on steps (b), (c) and (d), and assumes the feature mapping (i.e. learned through step (a)) is known and given.

et al., 2016; Sen et al., 2018; Hunziker et al., 2019), citizen science and crowdsourcing services (Sullivan et al., 2009; Nugent, 2018), or those targeting machine learning systems, such as model compression (Romero et al., 2014) and (understanding the vulnerability of) data poisoning attacks (Mei & Zhu, 2015; Zhu, 2018).

As illustrated in figure 1, a machine teaching framework assumes a computational model of the learner—either known or unknown to the teacher—which typically consists of two components: (a) a model for representing the learner's state (e.g., learner's current hypotheses), and (b) a model for the learning dynamics (e.g. parameters capturing learner's initial knowledge, learning rate, and learning behavior etc). When both models are known to the teacher, the teaching problem boils down to an optimal planning problem as in figure 1 (c). Upon receiving the teaching instructions, the learner makes an update according to its own intrinsic dynamics, and proceeds to the next knowledge state.

Classical theory of machine teaching often focuses on specific instantiations of such a framework. For example, assuming full access to the learner's dynamics and state representation, one may derive strong theoretical guarantees on the complexity of teaching (Goldman & Kearns, 1995; Zilles et al., 2011; Chen et al., 2018; Mansouri et al., 2019; Lessard et al., 2019; Tabibian et al., 2019; Hunziker et al., 2020). When the learner's representation is unknown but the learner's dynamics (e.g. learning algorithm) are given, it has been shown that the teacher can efficiently find a set of teaching examples with strong approximation guarantees in finding the optimal set (Dasgupta et al., 2019) or convergence guarantees (Liu et al., 2018) in teaching the target concept. When inspected under the teaching framework

¹Department of Computer Science, University of Chicago ²Department of Computer Science, University of Regina ³Max Planck Institute for Software Systems. Correspondence to: Chaoqi Wang <chaoqi@uchicago.edu>.

	representation			
model	known	unknown		
	white-box	"black-box learner"		
known	(Goldman & Kearns, 1995; Zilles et al., 2011; Chen et al., 2018)	(Dasgupta et al., 2019; Liu et al., 2018)		
	(Mansouri et al., 2019; Lessard et al., 2019; Tabibian et al., 2019; Hunziker et al., 2020)			
unknown	black-box MDP learner			
	(this work)	-		

Table 1: An overview of different teaching settings and difference between our work and the existing literature

in figure 1, these learner models can be viewed as special cases of *Markov learners*, where the goal is to reach the target state (e.g. the concept being taught). We list a few representative types of learners in section 2. Note that these existing results are often focused heuristic learner model or representations. Furthermore, these learner models can be classified under the category of non-epiphany learners (Chen & Krajbich, 2017; Dufwenberg et al., 2010), a class that may not always be suitable for modeling human behavior. Such assumptions and strictions hinder their applicability to solving practical problems, where the learner model is often a complicated blackbox (e.g. inferred from historic student data (Corbett & Anderson, 1994; Yudelson et al., 2013; Piech et al., 2015; Settles & Meeder, 2016; Sen et al., 2018; Hunziker et al., 2019)).

Our contributions. In this paper, we set forward a generic teaching framework capable of capturing the complex learner dynamics in the real-world teaching applications. In particular, we study machine teaching under a generic blackbox setting, where the learner's dynamics are modeled by a Markov decision process (MDP) with unknown parameters, covering both epiphany and non-epiphany learners. To provide a theoretical understanding, we derive the teaching cost under the assumption that the learning dynamics can be approximated by a linear function of the learner's state and the teaching instruction. Our contributions are:

- We introduce a generic machine teaching model that unifies many existing (heuristic) learner models, encompassing both epiphany and non-epiphany learners. This model allows us to learn the learner's model from data, providing a more versatile approach to machine teaching.
- Under our framework, we provide rigorous analyses on the teaching costs for various teaching scenarios. When the learning dynamics is linear, we show that the teaching costs grows at most polynomially in the optimal teaching cost and feature dimension d; when the dynamics is nonlinear, we show that teaching is not always feasible, and provide teachability conditions such that the teaching cost becomes linear (ignoring log factors) in the optimal cost.
- Complementing our theoretical results, we conduct experiments on a toy example to demonstrate the effectiveness of our proposed algorithm, and provide guidelines for setting its hyperparameters.

2. Related Work

Epiphany learning. The concept of Epiphany Learning (EL) has been rigorously studied in behavioral science (Chen & Krajbich, 2017; Dufwenberg et al., 2010). EL denotes a phenomenon where a learning agent (for instance, humans) experiences an abrupt enlightenment or comprehension regarding a specific subject matter. In the context of educational research, EL manifests when students achieve an insightful moment of comprehension or forge a substantial link between concepts, resulting in profound understanding of a topic or problem. Conversely, Non-Epiphany Learning implies scenarios in which such transformative moments do not transpire. Such epiphany/ non-epiphany learners naturally fit into the Markovian framework considered in this paper. We use MDP learner as a computational model to capture these learners, and subsequently provide an in-depth analysis of the teaching performance.

Markov learners in machine teaching Various learner models studied in the machine teaching literature can be viewed as Markov learners. Preference-based model for version space learners (Chen et al., 2018; Mansouri et al., 2019): Hereby, the state space corresponds to the version space and the learner's current hypothesis, action corresponds to teaching example, and transition probability is specified by the preference function σ . Gradient-based learner (Liu et al., 2017; 2018): Gradient-based learners are the major workhorse of machine learning. For them, the state space is \mathbb{R}^d , and the teaching set is the entire training set. Given the teaching example, the learner will be updated by the gradient descent rule, which is deterministic. In Liu et al. (2017; 2018), they study both the whitebox setting and the black-box setting. For the black-box setting, they assume that the learner's state is unknown but the transition function is known. Skill-based learner (Bower, 1961; Corbett & Anderson, 1994; Whitehill & Movellan, 2017): In the simplest form, the state space corresponds to d independent skills; at time step t, when an exercise x_t is presented, the skill associated with x_t can jump from 0 to 1 state with some probability (specified in table 2, notation adapted from Whitehill & Movellan (2017)). Memory-based learner (Hunziker et al., 2019): Classical computational models of human memory, such as the Half-Life Regression model (HLR), have been used in machine teaching to model

Type of the learner	States	Actions	Transition function
Preference-based version space (Mansouri et al., 2019)	$oldsymbol{h}_t \in 2^{\mathcal{H}} imes \mathcal{H}$	$\boldsymbol{x}_t \in \mathcal{X}$	$h_{t+1} \leftarrow \sigma(h_t)$
Gradient-based (Liu et al., 2017)	$oldsymbol{h}_t \in \mathbb{R}^d$	$\boldsymbol{x}_t \in \mathcal{X}$	$\boldsymbol{h}_{t+1} \leftarrow \boldsymbol{h}_t - \eta \cdot \nabla_{\boldsymbol{w}} \ell(\boldsymbol{h}_t, \boldsymbol{x}_t)$
Skill-based (Whitehill & Movellan, 2017)	$h_t \in [0, 1]^d$	$\boldsymbol{x}_t \in \mathcal{X}$	$oldsymbol{h}_{t+1} \propto oldsymbol{h}_t \odot g(oldsymbol{x}_t)^{lpha} \odot (1-g(oldsymbol{x}_t))^{(1-lpha)}$
Memory-based (Hunziker et al., 2019)	$oldsymbol{h}_t \in \mathbb{R}^2$	$\boldsymbol{x}_t \in \{0,1\}$	$oldsymbol{h}_{t+1} \leftarrow oldsymbol{h}_t \cdot ext{HL}(oldsymbol{x}_t)$

Table 2: Examples of existing sequential learner models that have Markov property.

the long term learning behavior of a human subject. State corresponds to the retention level and the forgetting rate, and transition is specified by the half-life dynamics HL. Nevertheless, these machine teaching models all assume the learner's model is known, and are designed in a heuristic way. In contrast, our work focus on proposing a generic framework that can capture these heuristic models and allow learning the learner's model from data.

Reinforcement teaching. Our work is also closely related to the reinforcement learning literature (Jaksch et al., 2010; Jin et al., 2020; Zhou et al., 2021; Ouyang et al., 2019; Min et al., 2021). In particular, our algorithm design is built upon the least-square regression algorithm for estimating the parameter of the dynamics function, and the extended value iteration (EVI) (Jaksch et al., 2010) for generating the teaching policy. These two sub-algorithms are commonly used as backbones in the algorithm design. In comparison, our work focus on the non-episodic setting and take the initialization into consideration, which better fits in the machine teaching problem. Another related line of works is teaching with reinforcement learning policy (Wu et al., 2018; Fan et al., 2018; Florensa et al., 2018; Omidshafiei et al., 2019). However, all of these works focus on improving the training efficiency of neural networks, i.e., whitebox learners. Their major contributions are developing better state representation and reward shaping functions based on different heuristics, which can serve as a complement to our work, i.e., the step (a) in figure 1.

3. Problem Formulation

We deal with the black-box setting, where the teacher initially has no knowledge of the learning dynamics of the learner, i.e., how the learner updates its knowledge state upon receiving the teacher instruction. We assume that the teacher can observe the learner's state directly, and also knows the cost function¹. The goal of the teacher is to help the learner reach some target knowledge state with minimal cost. To assist the learner, the teacher will not only provide informative teaching instructions to the learner but also needs to learn about the learner's dynamics.

Notations. Before we proceed, we first introduce some notation. We use \mathcal{H} to represent the set of all possible knowledge states of learners, h_0 denotes the initial knowledge state, and h_t is the learner's knowledge state at iteration t. At each iteration t, the teacher can choose one teaching instruction x_t from the teaching set \mathcal{X} . The learner's knowledge state will be updated upon receiving the teaching instruction. The teacher's goal is to help the learner transit to the target knowledge state h^* with minimal cost. Throughout the entire paper, we use C_* to denote the tightest upper bound on the expected teaching cost of the optimal teaching policy by starting from any initial state.

3.1. Parametric Markov Learners

We model the learner as a Markov learner, which is able to cover a broad class of learners considered in the literature (Gao et al., 2016; Whitehill & Movellan, 2017; Liu et al., 2017; Hunziker et al., 2019; Mansouri et al., 2019). Specifically, for any given learner, it starts from some initial knowledge state h_0 , which represents its current knowledge state. For each iteration t, when the learner receives the teaching instruction x_t , it updates its knowledge based on its transition probability,

$$\boldsymbol{h}_{t+1} \sim \mathbb{P}_{\boldsymbol{\theta}^{\ell}}[\boldsymbol{h}_{t+1}|\boldsymbol{h}_{t}, \boldsymbol{x}_{t}],$$
 (1)

where θ^ℓ refers to the parameters that define the transition probability or learning dynamics of the learner. Different learners may have different θ^ℓ . The transition probability induces a preference over the next knowledge states for the learner, which captures the learning dynamics of the learner. Intuitively, if a learner is smart enough, then the learner will assign a higher transition probability to states that are close to the target state h^\star upon receiving the teaching instructions. In contrast, sometimes, a learner may not be able to understand advanced teaching instructions before it reaches some knowledge state. To model such scenarios, the learner may assign a very high probability to remain at the current knowledge state when receiving obscure teaching instructions.

3.2. The Teacher's Objective

The teacher's goal is to help the learner learn as fast as possible, i.e., minimizing the cost of steering the learner to reach the target knowledge state h^* . In order to teach,

¹In practice, the teacher can probe the learner's knowledge state by quizzing the learner. The cost could be the price of the teaching instruction.

there are two tasks that the teacher needs to solve, namely estimating θ^{ℓ} and generating the teaching instruction. The entire problem can be formulated as follows, where $c(\cdot, \cdot)$ is the cost function.

$$\min_{\boldsymbol{x}_{1:T} \in \mathcal{X}^T, T \in \mathbb{Z}_+} \quad \sum_{t=1}^T c(\boldsymbol{h}_t, \boldsymbol{x}_t)$$
 (2)

s.t.
$$h_{T+1} = h^*$$
 and $h_{t+1} \sim \mathbb{P}_{\theta^{\ell}}(h|h_t, x_t)$.

If the teacher knows the true parameters, then the above problem becomes a (stochastic) planning problem. In this work, we assume that the teacher only knows the parametric form of the learner's transition function, and it doesn't know the true parameters of the learner. This introduces an extra challenge in solving the above problem, but it also makes the problem formulation more general.

4. Teaching Blackbox Markov Learners: Algorithm and Analysis

In this section, we present an algorithm for teaching blackbox Markov learners (including epiphany learners and nonepiphany learners), which takes the initialization into account. We first introduce the assumptions that the subsequent sections rely on in subsection 4.1. Then, we conduct a rigorous analysis for upper bounding the teaching cost under different teaching scenarios, including 1) the Markov learner is linear and teachable; 2) the Markov learner is nonlinear and teachable.

4.1. Preliminaries and Background

Teaching Markov learners can be captured by an MDP M := $\{\mathcal{H}, \mathcal{X}, \mathbb{P}, c, \mathbf{h}_0, \mathbf{h}^*, \gamma\}$, where $c: \mathcal{H} \times \mathcal{X} \to \mathbb{R}_+$ is the cost function and h^\star is the target knowledge state. For any $(h, x, h') \in \mathcal{H} \times \mathcal{X} \times \mathcal{H}$, $\mathbb{P}_{\theta^{\ell}}(h'|h, x)$ denotes the probability of transiting to knowledge state h' given the teaching instruction x under h. To be noted, when the learner reaches the target knowledge state, the cost will be zero for all the teaching instructions, i.e., $c(\mathbf{h}^{\star}, \mathbf{x}) =$ $0, \ \forall \ \boldsymbol{x} \in \mathcal{X}, \text{ and } \mathbb{P}(\boldsymbol{h}^{\star} | \boldsymbol{h}^{\star}, \boldsymbol{x}) = 1, \text{ which means the target}$ knowledge state is an absorbing state. $\gamma \in (0,1]$ is the cost discounting factor. In the teaching context, $1 - \gamma$ is the probability of the learner transiting to the target knowledge state from any other state, i.e., the probability of epiphany learning (Dufwenberg et al., 2010; Chen & Krajbich, 2017).

Definition 1 (Proper Policy) A stationary policy π is proper if, given any initial state, the probability of reaching the goal state g within a finite number of steps, when following π , is strictly positive.

Let us denote by $\Pi_{proper}(M)$ the set of stationary polices of the underlying MDP M such that for any policy

 $\pi \in \Pi_{proper}(M)$, the expected time that it takes to reach the target knowledge state h^{\star} from any initial knowledge state h is finite. In the teaching context, the existence of proper polices for a learner means that there is a way to teach him/her the target knowledge state h^* . In our analvsis, we will assume that the Markov learner is linear and teachable under some known and pregiven feature mapping $\phi: \mathcal{H} \times \mathcal{X} \times \mathcal{H} \to \mathbb{R}^d$. We summarize the essential idea in the following assumption.

Assumption 1 (Teachable Linear Markov Learners)

 $M := \{\mathcal{H}, \mathcal{X}, \mathbb{P}_{\theta^{\ell}}, c, h_0, h^{\star}, \gamma\}$ is a teachable linear Markov learner, if it satisfies

- *Linearity*: Given a known feature mapping ϕ , there exists an unknown parameter $\boldsymbol{\theta}^{\ell} \in \mathbb{R}^d (\|\boldsymbol{\theta}^{\ell}\|_2^2 \leq d)$ such that $\mathbb{P}_{\boldsymbol{\theta}^{\ell}}(\boldsymbol{h}'|\boldsymbol{h},\boldsymbol{x}) = \langle \boldsymbol{\phi}(\boldsymbol{h}'|\boldsymbol{h},\boldsymbol{x}), \boldsymbol{\theta}^{\ell} \rangle, \forall (\boldsymbol{h},\boldsymbol{x},\boldsymbol{h}') \in$ $\mathcal{H} \times \mathcal{X} \times \mathcal{H}$.
- Teachable: There exists at least one proper policy, i.e., $\Pi_{proper}(M) \neq \emptyset$.

Furthermore, for any bounded value function V: $\mathcal{H} \rightarrow [0,C] \text{ with } C_{\star} \leq C, \|\phi_{V}(\boldsymbol{h},\boldsymbol{x})\|_{2} \leq \sqrt{d}C$ holds for any $(h,x) \in \mathcal{H} \times \mathcal{X}$, where $\phi_V(h,x) =$ $\sum_{h'} \phi(h'|h,x)V(h')$.

For any value function $V:\mathcal{H}\to\mathbb{R}_+,$ we define $\mathbb{P}V(\pmb{h},\pmb{x}) = \sum_{\pmb{h}'} \mathbb{P}(\pmb{h}'|\pmb{h},\pmb{x}) V(\pmb{h}') \text{ for any } (\pmb{h},\pmb{x}) \in \mathcal{H} \times \mathcal{H}$ \mathcal{X} . Under the linear MDP assumption, we further have $\mathbb{P}_{\boldsymbol{\theta}^{\ell}}V(\boldsymbol{h},\boldsymbol{x}) = \langle \boldsymbol{\phi}_V(\boldsymbol{h},\boldsymbol{x}), \boldsymbol{\theta}^{\ell} \rangle$. For convenience of notation, we further define the cost-to-go function for policy π under $M_{\theta\ell}$ as

$$V^{\pi}(\boldsymbol{h}|\boldsymbol{\theta}^{\ell}) \coloneqq \lim_{T \to +\infty} \mathbb{E}\left[\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \middle| \boldsymbol{h}_{0} = \boldsymbol{h}\right], \quad (3)$$

$$m{h}_{t+1} \sim \mathbb{P}_{m{ heta}^{\ell}}(m{h}|m{h}_t, m{x}_t) \text{ and } m{x}_t = \pi(m{h}_t).$$

Consequently, the Q-value function of policy π under M_{θ^ℓ}

$$Q^{\pi}(\boldsymbol{h}, \boldsymbol{x} | \boldsymbol{\theta}^{\ell}) := c(\boldsymbol{h}, \boldsymbol{x}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V^{\pi}(\boldsymbol{h}, \boldsymbol{x} | \boldsymbol{\theta}^{\ell}). \tag{4}$$

Subsequently, when there is no ambiguity, we use V(h) and Q(h, x) to simplify notation. Next, we introduce another assumption tailored to the teaching setting.

Assumption 2 (δ_0 -Closeness) The true parameter θ^{ℓ} is δ_0 close to the teacher's initial estimation θ_0 , i.e., $\|\theta^{\ell} - \theta_0\|_2 \le$ $\delta_0 \sqrt{d}$ with $0 \le \delta_0 \le 1$.

The above assumption is natural in the teaching setting. Without such an assumption, the teacher may need to interact with the learner for a large number of rounds before it can teach in an effective way, which is impractical for teaching resource-constrained learners, such as humans. In practice, to fulfil Assumption 2, we can first fit a transition

Algorithm 1 Blackbox Teaching Algorithm for Non-Epiphany and Epiphany Learners.

```
Require: Initial estimation \hat{\theta}_0 = \theta_0, iteration t = 0, EVI
        index t_0 = 0, k = 0, \Sigma_0 = \lambda I, \mu_0 = \lambda \theta_0 and discount
        factor \gamma (for epiphany learners).
  1: Q_0 \leftarrow EVI(\{\hat{\boldsymbol{\theta}} \in \mathbb{R}^d | \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_0\|_2 \leq \delta_0\}, \frac{1}{\lambda}, \frac{1}{\lambda})
  2: while h_t \neq h^* do
             Provide
                                      teaching instruction x_t
             \arg\min_{\boldsymbol{x}\in\mathcal{X}}Q_k(\boldsymbol{h}_t,\boldsymbol{x}); Receive c_t=c(\boldsymbol{h}_t,\boldsymbol{x}_t) and
             h_{t+1} \sim \mathbb{P}_{\boldsymbol{\theta}^{\ell}}(\cdot|\boldsymbol{h}_t, \boldsymbol{x}_t).
             oldsymbol{\Sigma}_t \leftarrow oldsymbol{\Sigma}_{t-1} + oldsymbol{\phi}_{V_k}(oldsymbol{h}_t, oldsymbol{x}_t) oldsymbol{\phi}_{V_k}(oldsymbol{h}_t, oldsymbol{x}_t)^{	op}
             \boldsymbol{\mu}_t \leftarrow \boldsymbol{\mu}_{t-1} + \boldsymbol{\phi}_{V_k}(\boldsymbol{h}_t, \boldsymbol{x}_t) V_k(\boldsymbol{h}_{t+1})
  5:
             if det(\Sigma_t) \geq 2 det(\Sigma_{t_k}) or t \geq 2t_k + \lambda then
  6:
  7:
                  k \leftarrow k + 1
  8:
                  9:
10:
11:
12:
             t \leftarrow t + 1
13: end while
```

function on the offline teacher-learner interaction data to serve as the initialization. For simplicity, we denote the associated MDP of a learner with parameter θ^{ℓ} as $M_{\theta^{\ell}}$ and the teacher's initial estimation on the parameter as θ_0 .

Lastly, we define two categories of Markov learners depending on their behaviors during learning, which can be precisely modelled by undiscounted MDP and discounted MDP, respectively. We call a Markov learners an *epiphany learner* if $\gamma < 1$ for its associated MDP. When the learner's associated MDP has $\gamma = 1$, we call it a *non-epiphany learner*. Epiphany learning (Dufwenberg et al., 2010; Chen & Krajbich, 2017) is an observed phenomenon in human learners, which is often depicted by a light bulb showing over a person's head in cartoons. In the context of machine teaching, $1 - \gamma$ can be intuitively understood as the lower bound of the probability of epiphany learning (i.e., transit to the goal knowledge state) at all the knowledge states.

4.2. Blackbox Teaching for Linear Markov Learners

We first consider the case where the Markov learner is linear and teachable (see Assumption 1). We first present an algorithm for solving the teaching problem, which takes the initialization into consideration. The entire algorithm is built upon solving a regularized least-squares regression (for computing $\hat{\theta}$), and extended value iteration (for generating the teaching policy). These two sub-algorithms are often used as backbones for algorithm design (Jaksch et al., 2010; Jin et al., 2020; Zhou et al., 2021; Ouyang et al., 2019; Min et al., 2021). In contrast, our algorithm 1) takes the initializa-

Algorithm 2 Extended Value Iteration: EVI(\mathcal{C}, ξ, ν)

```
Require: Confidence set C, error tolerance of valute itera-
       tion \xi, iteration i=0, cost discount factor \nu.
  1: Q^{(0)}(\cdot, \cdot) = 0
  2: \ Q(\cdot, \cdot) = 0
  3: V^{(0)}(\cdot) = 0
  4: V^{(-1)}(\cdot) = +\infty
  5: if \mathcal{C} \cap \mathcal{B} \neq \emptyset then
            6:
  7:
                 \begin{array}{l} \min_{\boldsymbol{\theta} \in \mathcal{C} \cap \mathcal{B}} \langle \boldsymbol{\theta}, \boldsymbol{\phi}_{V^{(i)}}(\cdot, \cdot) \rangle \\ V^{(i+1)}(\cdot) \leftarrow \min_{\boldsymbol{x} \in \mathcal{X}} Q^{(i+1)}(\cdot, \boldsymbol{x}) \end{array} 
  8:
  9:
                 i \leftarrow i + 1
10:
            end while
            Q(\cdot, \cdot) \leftarrow Q^{(i+1)}(\cdot, \cdot).
11:
12: end if
13: return Q(\cdot, \cdot)
```

tion θ_0 into account, which is crucial to teaching effectively; 2) and applies to both epiphany and non-epiphany learners. Intuitively, Algorithm 1 can be divided into two parts as described below.

Parameter learning. For parameter learning, once the updating criteria is satisfied, the teacher will update its estimation of the learner's parameter based on the interactions so far. Updating the estimation reduces to solving the following initialization-regularized least-squares problem:

$$\hat{\boldsymbol{\theta}}_{m} \leftarrow \operatorname*{arg\,min}_{\boldsymbol{\theta} \in \mathbb{R}^{d}} \sum_{t=0}^{m-1} \left[\langle \boldsymbol{\phi}_{V_{k(t)}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}), \boldsymbol{\theta} \rangle - V_{k(t)}(\boldsymbol{h}_{t+1}) \right]^{2} + \lambda \|\boldsymbol{\theta} - \boldsymbol{\theta}_{0}\|_{2}^{2}, \tag{5}$$

where k(t) is the index of the value function at iteration t, e.g., for $t_{j-1} \le t \le t_j - 1$, the index is k(t) = j - 1, and $V_j(h)$ is the j_{th} value function returned by the extended value iteration (EVI) algorithm (Jaksch et al., 2010). The above problem has a closed-form solution $\hat{\theta}_m = \Sigma_m^{-1} \mu_m$, where (see also Lines 4&5 in Algorithm 1),

$$oldsymbol{\Sigma}_m = \lambda oldsymbol{I} + \sum_{t=0}^{m-1} oldsymbol{\phi}_{V_{k(t)}}(oldsymbol{h}_t, oldsymbol{x}_t) oldsymbol{\phi}_{V_{k(t)}}(oldsymbol{h}_t, oldsymbol{x}_t)^{ op}, \quad (6)$$

$$\mu_m = \lambda \theta_0 + \sum_{t=0}^{m-1} \phi_{V_{k(t)}}(\mathbf{h}_t, \mathbf{x}_t) V_{k(t)}(\mathbf{h}_{t+1}).$$
 (7)

The value of λ indicates our confidence on the optimality of the initialization θ_0 . When the initialization is very likely close to the true parameter θ^{ℓ} , we should set a large λ , otherwise we should set a small λ . In addition, λ also affects the updating frequency of the parameter, which is triggered by two criteria, namely 1) the log-determinant of Σ_t ; and 2) the number of iterations (see Line 7 in Algorithm 1). When

 λ is larger, the parameter will be updated less frequently, since we trust our current estimate more. To be noted, in our analysis, we always assume $\lambda \geq 1^2$.

Teaching. During the teaching phase, the teacher's policy is induced by the Q-value function returned by EVI (see Algorithm 2). After the teacher's teaching instruction, the teacher will receive a cost incurred by the teaching instruction, and also observe the learner's latest knowledge state,

$$\boldsymbol{x}_{t} = \arg\min_{\boldsymbol{x} \in \mathcal{X}} Q_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}), \tag{8}$$

$$c_t = c(\boldsymbol{h}_t, \boldsymbol{x}_t), \quad \boldsymbol{h}_{t+1} \sim \mathbb{P}_{\boldsymbol{\theta}^{\ell}}(\boldsymbol{h}|\boldsymbol{h}_t, \boldsymbol{x}_t),$$
 (9)

In detail, the EVI algorithm takes the confidence set \mathcal{C}_t (see Lemma 1 for $t \geq 1$), the error tolerance of the value iteration ξ and the cost discounting factor ν as input. The confidence set \mathcal{C}_t is an ellipsoid centered at the current estimation $\hat{\theta}_t$. With high probability, the true parameter θ^{ℓ} lies in the intersection of \mathcal{B} and \mathcal{C}_t , where \mathcal{B} is defined as

$$\mathcal{B} \coloneqq \{ oldsymbol{ heta} \in \mathbb{R}^d | \langle oldsymbol{\phi}(\cdot | oldsymbol{h}, oldsymbol{x}), oldsymbol{ heta}
angle \in \Delta^d, \ orall (oldsymbol{h}, oldsymbol{x}) \in \mathcal{H} imes \mathcal{X} \}.$$

The error tolerance parameter is chosen to be $\xi=1/(\lambda t)$. Intuitively, the error tolerance will be smaller when we 1) collect more data (i.e., t becomes large); and 2) start from a better initialization (i.e., δ_0 is smaller). For the cost discount factor ν , we set it to be $1-1/(\lambda t)$, when the underlying MDP of the Markov learner is undiscounted (i.e., non-epiphany learners). By doing so, the cost discount factor ν will become closer to 1 as the teaching continues, which helps us avoid a teaching cost that is linear in T (i.e, the total number of teaching instructions) and also ensures the convergence of EVI. When the learner's underlying MDP is discounted (i.e., epiphany learners), we will set $\nu = \gamma$ to be a constant. Intuitively, the cost discount factor ν captures the probability of epiphany learning.

Overall, the EVI algorithm adapts the standard value iteration algorithm to incorporate the optimism-in-the-face-of-uncertainty (OFU) principle (see Line 7 in Algorithm 2) proposed by Abbasi-Yadkori et al. (2011), which has been demonstrated to be effective in online learning setting.

4.3. Theoretical Analysis for the Linear Case

In this section, we analyze the cost upper bounds of using Algorithm 1 for teaching both non-epiphany learners and epiphany learners. The core of the algorithm is to build the confidence set that contains the true parameter θ^ℓ , which balances exploration (parameter learning) and exploitation (teaching) automatically. In general, the smaller the confidence set that we can construct, the lower the cost. In the following, we present Lemma 1, which provides a confidence set containing θ^ℓ with high probability.

Lemma 1 Under Assumptions 1 and 2, for any $t \ge 1$, with probability at least $1 - \delta$, we have that the true parameter θ^{ℓ} lies in

$$C_t = \left\{ \boldsymbol{\theta} \in \mathbb{R}^d \,\middle|\, \|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}\|_{\boldsymbol{\Sigma}_t} \le r_t \right\},\tag{10}$$

where
$$r_t := C\sqrt{d\log((4(t^2 + t^3C^2/\lambda))/\delta)} + \sqrt{\lambda}\delta_0$$
.

The confidence set \mathcal{C}_t is centered at the current estimation $\hat{\theta}_t$. Its radius is computed based on the iteration t, feature dimension d, regularization parameter λ , the upper bound of the optimal cost C, and the upper bound on the distance between θ_0 and θ^ℓ , i.e., δ_0 . The confidence set will become smaller as δ_0 decreases, indicating that a good initialization is desired. Theorem 1 provides an upper bound on the cost of teaching non-epiphany learners using Algorithm 1.

Theorem 1 Under Assumptions 1 and 2, if the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(C_\star)$, $\lambda = 1/\delta_0^2$, and the cost function is bounded from below by c_{\min} for all non-goal knowledge states $(\mathcal{H} \setminus \{h^\star\})$ and teaching instruction (\mathcal{X}) pairs, then with probability at least $1 - 2\delta$, the teaching cost of Algorithm 1 for non-epiphany learners (i.e., $\gamma = 1$) is upper bounded by

$$\mathcal{O}\left(\left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}d\delta_{0}}{\delta c_{\min}}\right)}\right) \cdot \log^{1.5}\left(\frac{C_{\star}d}{c_{\min}\delta}\right) \cdot \frac{C_{\star}^{2}d}{c_{\min}}\right). \tag{11}$$

The cost upper bound in Theorem 1 has a polynomial dependency on the expected cost of the optimal policy, C_{\star} . It's worth noting that when $\delta_0 \to 0$, the purple term inside the parentheses of Equation 11 will vanish leaving only the constant term 1. The constant term 1 is due to the stochasticity in the transition of the learner's knowledge states, which is independent of the teaching algorithm used.

Next, we consider the case where the learner is an *epiphany learner*. Intuitively, epiphany learning can be interpreted as adding a shortcut from the current knowledge state to the target knowledge state in the underlying MDP, which is equivalent to the discounted MDP case. The following theorem provides an upper bound on the cost of teaching epiphany learners.

Theorem 2 Under Assumptions 1 and 2, if the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(C_\star)$, $\lambda = 1/\delta_0^2$, then with probability at least $1 - 3\delta$, the total cost incurred by running Algorithm 1 for epiphany learners with $\gamma < 1$, is upper bounded by

$$\mathcal{O}\left(C_{\star} \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}^{2} \delta_{0}^{2} \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C_{\star} \log \delta}{\delta \log \gamma}\right)}\right). \tag{12}$$

²This is because we found that $\lambda \propto 1/\delta_0^2$ is a good choice in practice, and $\delta_0 \leq 1$ by Assumption 2

Compared with Theorem 1 for non-epiphany learners, the upper bound of the teaching cost for epiphany learners is linear (ignoring the log factors) in the expected teaching cost of the optimal policy C_{\star} and the feature dimension d. Moreover, the dependency on d will vanish when $\delta_0 \to 0$ as well. In addition, Theorem 2 does not require the cost function to be bounded from below.

4.4. Theoretical Analysis for the Non-linear Case

In the previous section, we presented the theoretical analysis for both non-epiphany and epiphany learners when their learning dynamics are linear. One natural follow-up question is: what would happen if the learner's dynamics is non-linear, i.e., the linear model is misspecified? To study this problem, we consider the case where teaching the learner can be *approximately* modelled as a linear MDP. This idea is captured in the following assumption.

Definition 2 (ε-Approximate Teachable Markov Learners)

For any $\epsilon \in (0,1]$, a MDP $M = (\mathcal{H}, \mathcal{X}, \mathbb{P}, c, \mathbf{h}_0, \mathbf{h}^*, \gamma)$ is an ϵ -approximate teachable MDP with a feature map ϕ , if there exists a unknown teachable linear MDP M_{θ^*} such that for any $(\mathbf{h}, \mathbf{x}) \in \mathcal{H} \times \mathcal{X}$, we have $\|\mathbb{P}(\cdot|\mathbf{h}, \mathbf{x}) - \langle \phi(\cdot|\mathbf{h}, \mathbf{x}), \theta^* \rangle\|_{TV} \leq \epsilon$, where TV denotes the total variation distance.

By definition, the learner is an ϵ -approximate teachable Markov learner if the learning dynamics function of the learner is close to a linear transition function under the given feature mapping ϕ . We measure the closeness between the dynamics functions by the total variation distance.

In general, the algorithm designed for the linear case will fail when the transition function is non-linear. Specifically, for non-epiphany learners, the teaching cost can be unbounded even for a small model misspecification level ϵ . To illustrate this, we present an informal example (see Figure 2), where the teaching policy induced by the closest linear MDP to the learner's MDP will incur an infinite teaching cost. The intuition behind such counterexamples is that the teaching policy induced by the misspecified MDP will get trapped in a circle of the true MDP. Fortunately, for epiphany learners, the teaching cost of Algorithm 1 can still be bounded well, and it is robust to small misspecification levels. The results are stated in the following theorem.

Theorem 3 For ϵ -approximate teachable epiphany learners as defined in Definition 2, if $\|\theta_0 - \theta^*\|_2 \le \delta_0$, the cost function is bounded from above by c_{max} , the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(\epsilon \gamma c_{max}/(1-\gamma)^2 + C_*)$, and if $\lambda = 1/\delta_0^2$, then with probability at least $1 - 3\delta$, the teaching cost incurred by

running Algorithm 1 is upper bounded by

$$\mathcal{O}\left(C \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C^2 \delta_0^2 \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C \log \delta}{\delta \log \gamma}\right)} + \frac{\epsilon \log \delta}{\log \gamma}C\right).$$
(13)

In contrast to Theorem 2, the major difference is that there is one extra cost term in Theorem 3 due to the intrinsic bias of the linear approximation. When ϵ is sufficiently small, those terms with coefficient ϵ can be ignored safely, which gives us the following proposition.

Proposition 1 Under the same assumptions as Theorem 3, if $\epsilon = \mathcal{O}\left(C_{\star}(1-\gamma)^2/(\gamma c_{max})\right)$ then with probability at least $1-3\delta$, the total cost incurred by running Algorithm 1 is upper bounded by

$$\mathcal{O}\left(C_{\star} \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}^{2} \delta_{0}^{2} \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C_{\star} \log \delta}{\delta \log \gamma}\right)} + \frac{\epsilon \log \delta}{\log \gamma}C_{\star}\right).$$

$$(14)$$

Hence, as indicated by Theorem 3 and Proposition 1, our algorithm can still attain good theoretical guarantees when the misspecification level is low.

5. Numerical Experiments

Our experiments serve as complements to our theories, which also provides a sanity check for our algorithm. Firstly, we present experiments to examine the performance of the proposed algorithm. In addition, we also evaluate how the choice of λ affects the empirical teaching cost, as λ plays a critical role in our algorithm design.

5.1. Experimental Setup

Knowledge states and teaching instructions. We sample 100 weights $\{\boldsymbol{h}_i\}_{i=1}^{100}$ uniformly at random from $[-3,3]^d$ to simulate different knowledge states, each of which corresponds to a linear regressor. We then pick one of the weights to represent the target knowledge state, denoted as \boldsymbol{h}^* . To generate the teaching instructions, we first sample 20 points $\{\boldsymbol{z}_i\}_{i=1}^{20}$ from a normal distribution $\mathcal{N}(\boldsymbol{0},\boldsymbol{I})$, and their corresponding labels are generated by $y_i = \langle \boldsymbol{h}^*, \boldsymbol{z}_i \rangle + \zeta$, where $\zeta \sim \mathcal{N}(0,1)$ is the observation noise. By $\{\boldsymbol{x}_i\}_{i=1}^{20}$ we denote the set of teaching instructions, where $\boldsymbol{x}_i = (\boldsymbol{z}_i, y_i)$.

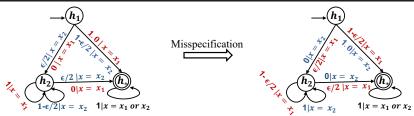


Figure 2: An illustration of the failure case under the misspecified setting. The MDP consists of 2 teaching actions $\mathcal{X} = \{x_1, x_2\}$ and 3 states $\mathcal{H} = \{h_0, h_2, h^*\}$, and the misspecification level is ϵ . For the teaching policy induced by the misspecified MDP (right), the learner can get stuck at the state h_2 with probability $\epsilon/2$.

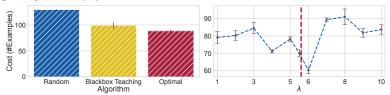


Figure 3: Left: A comparison between random teaching policy, blackbox teaching policy and the optimal teaching policy in terms of the mean of the averaged teaching cost for 99 initial states; Right: Effect of different values of λ on the teaching cost (computed with the first 10 states to save computation time).

Feature representation. We consider the feature representation for each triplet (h, x, h') to be

$$\phi(\boldsymbol{h}'|\boldsymbol{h},\boldsymbol{x}) = \begin{bmatrix} 1/\left(Z_{(\boldsymbol{h},\boldsymbol{x})}^{(1)} \cdot \|\boldsymbol{h}' - \boldsymbol{h} + \eta \nabla_{\boldsymbol{h}} \ell(\boldsymbol{h},\boldsymbol{x})\|_{2}\right) \\ 1/\left(Z_{(\boldsymbol{h},\boldsymbol{x})}^{(2)} \cdot \|\nabla_{\boldsymbol{h}} \ell(\boldsymbol{h},\boldsymbol{x})\|_{2}\right) \end{bmatrix}$$

where η is the learning rate, and $Z^{(i)}_{(\boldsymbol{h},\boldsymbol{x})}$ is the normalizing constant for the i^{th} dimension of the feature representation $\phi(\cdot|\boldsymbol{h},\boldsymbol{x})$. The normalizing constants are used to ensure that $\sum_{\boldsymbol{h}'\in\mathcal{H}}\phi(\boldsymbol{h}'|\boldsymbol{h},\boldsymbol{x})=(1,1)$. Therefore, all the feasible θ that forms a probabilistic distribution lies in a 1-d simplex. Intuitively, the first dimension indicates that the learner is more likely to transit to those knowledge states that align well with the updated knowledge state, i.e., $\boldsymbol{h}-\eta\nabla_{\boldsymbol{h}}\ell(\boldsymbol{h},\boldsymbol{x})$, whereas the second dimension implies that the learner's knowledge state transition will become more random if the teaching instruction is more difficult, which is measured by the gradient norm $\|\nabla_{\boldsymbol{h}}\ell(\boldsymbol{h},\boldsymbol{x})\|_2$.

5.2. Empirical Results

Comparing with baselines. We first evaluate the empirical performance of Algorithm 1 under the above experimental setup. Specifically, we set the learning rate $\eta=1$, and compare it with the random teaching policy and the optimal policy. We compute the mean of the averaged teaching cost of starting from each non-goal knowledge state. The averaged teaching cost is computed with 50 random seeds. The results are presented in figure 3 left. We see that the blackbox teaching algorithm outperforms the random teaching policy but underperforms the optimal teaching policy.

How to set λ ? The initialization plays an important role in our algorithm design and theoretical analysis. Our theoretical analysis has demonstrated the impact of the initialization on the teaching cost. However, given the initialization, it is still unclear how to set the right regularization parameter λ .

We conjecture that the 'optimal' λ should be around $1/\delta_0^2$, which is also adopted in our theoretical analysis. To verify this idea, we study how the choice of λ affects the teaching cost. Under the same setting as above, we vary the value of λ in $\{1.0, 1.5, ..., 10.0\}$. To save computation time, we adopt the first 10 states to serve as the initial state and repeat the previous experiments. The results are reported in the right plot of figure 3. The red dashed line corresponds to the line of $x=1/\delta_0^2$ with $\delta_0^2=0.18$. Based on the empirical results, we can observe that the best choice of λ is 6, which is close to $1/\delta_0^2$. In addition, if we set λ too large or too small, the teaching cost will increase accordingly.

In summary, our experimental results highlight that modelling the learner's learning dynamics is crucial to achieve a low teaching cost. Furthermore, given the initialization, setting $\lambda=1/\delta_0^2$ is a reasonable choice for obtaining good empirical performance.

6. Discussion and Conclusion

In this paper, we investigate a generic framework for machine teaching, under which the learner's dynamics can be represented as an MDP with unknown, learnable parameters. To solve the teaching problem, we introduce an algorithm that accommodates both epiphany and non-epiphany learners, thus bridging a significant gap in the current literature. Moreover, we furnish a rigorous analysis of the teaching costs associated with these two types of learners under disparate settings. Complementing our theoretical insights, we conduct empirical research to demonstrate the efficiency of our proposed algorithm and provide a guideline for setting hyperparameters. It is our aspiration that this work will stimulate future research in proposing more nuanced assumptions about the structure of the learner's MDP and more efficient algorithms for machine teaching.

References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Bower, G. H. Application of a model to paired-associate learning. *Psychometrika*, 26(3):255–280, 1961.
- Chen, W. J. and Krajbich, I. Computational modeling of epiphany learning. *Proceedings of the National Academy* of Sciences, 114(18):4637–4642, 2017.
- Chen, Y., Singla, A., Mac Aodha, O., Perona, P., and Yue, Y. Understanding the role of adaptivity in machine teaching: The case of version space learners. *Advances in Neural Information Processing Systems* 32, 2018.
- Corbett, A. T. and Anderson, J. R. Knowledge tracing: Modeling the acquisition of procedural knowledge. *User modeling and user-adapted interaction*, 1994.
- Csáji, B. C. and Monostori, L. Value function based reinforcement learning in changing markovian environments. *Journal of Machine Learning Research*, 9(8), 2008.
- Dasgupta, S., Hsu, D., Poulis, S., and Zhu, X. Teaching a black-box learner. In *International Conference on Machine Learning*, pp. 1547–1555. PMLR, 2019.
- Dufwenberg, M., Sundaram, R., and Butler, D. Epiphany in the game of 21. *Journal of Economic Behavior & Organization*, 75(2):132–143, 2010.
- Fan, Y., Tian, F., Qin, T., Li, X.-Y., and Liu, T.-Y. Learning to teach. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=HJewuJWCZ.
- Florensa, C., Held, D., Geng, X., and Abbeel, P. Automatic goal generation for reinforcement learning agents. In *International conference on machine learning*, pp. 1515–1528. PMLR, 2018.
- Gao, Z., Ries, C., Simon, H., and Zilles, S. Preferencebased teaching. In *Conference on Learning Theory*, pp. 971–997. PMLR, 2016.
- Goldman, S. A. and Kearns, M. J. On the complexity of teaching. *Journal of Computer and System Sciences*, 50 (1):20–31, 1995.
- Hunziker, A., Chen, Y., Mac Aodha, O., Rodriguez, M. G., Krause, A., Perona, P., Yue, Y., and Singla, A. Teaching multiple concepts to a forgetful learner. In *NeurIPS*, 2019.
- Hunziker, A., Chen, Y., Mac Aodha, O., Gomez Rodriguez,M., Krause, A., Perona, P., Yue, Y., and Singla, A. Teaching multiple concepts to a forgetful learner. *Advances in*

- Neural Information Processing Systems 32, 6:4025–4036, 2020.
- Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010. URL http://jmlr.org/papers/v11/jaksch10a.html.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143. PMLR, 2020.
- Lessard, L., Zhang, X., and Zhu, X. An optimal control approach to sequential machine teaching. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2495–2503. PMLR, 2019.
- Liu, W., Dai, B., Humayun, A., Tay, C., Yu, C., Smith, L. B., Rehg, J. M., and Song, L. Iterative machine teaching. In *International Conference on Machine Learning*, pp. 2149–2158. PMLR, 2017.
- Liu, W., Dai, B., Li, X., Liu, Z., Rehg, J., and Song, L. Towards black-box iterative machine teaching. In *International Conference on Machine Learning*, pp. 3141–3149. PMLR, 2018.
- Mansouri, F., Chen, Y., Vartanian, A., Zhu, X., and Singla, A. Preference-based batch and sequential teaching: Towards a unified view of models. *Advances in Neural Information Processing Systems* 32, 2019.
- Mei, S. and Zhu, X. Using machine teaching to identify optimal training-set attacks on machine learners. In *AAAI*, pp. 2871–2877, 2015.
- Min, Y., He, J., Wang, T., and Gu, Q. Learning stochastic shortest path with linear function approximation. *arXiv* preprint arXiv:2110.12727, 2021.
- Nugent, J. iNaturalist: citizen science for 21st-century naturalists. *Science Scope*, 41(7):12, 2018.
- Omidshafiei, S., Kim, D.-K., Liu, M., Tesauro, G., Riemer, M., Amato, C., Campbell, M., and How, J. P. Learning to teach in cooperative multiagent reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 6128–6136, 2019.
- Ouyang, Y., Gagrani, M., and Jain, R. Posterior sampling-based reinforcement learning for control of unknown linear systems. *IEEE Transactions on Automatic Control*, 65(8):3600–3607, 2019.
- Piech, C., Bassen, J., Huang, J., Ganguli, S., Sahami, M., Guibas, L. J., and Sohl-Dickstein, J. Deep knowledge tracing. In NIPS, 2015.

- Rafferty, A. N., Brunskill, E., Griffiths, T. L., and Shafto, P. Faster teaching via pomdp planning. *Cognitive science*, 40(6):1290–1332, 2016.
- Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550, 2014.
- Sen, A., Patel, P., Rau, M. A., Mason, B., Nowak, R., Rogers, T. T., and Zhu, X. Machine beats human at sequencing visuals for perceptual-fluency practice. In *EDM*, 2018.
- Settles, B. and Meeder, B. A trainable spaced repetition model for language learning. In *ACL*, pp. 1848–1858, 2016.
- Sullivan, B. L., Wood, C. L., Iliff, M. J., Bonney, R. E., Fink, D., and Kelling, S. eBird: A citizen-based bird observation network in the biological sciences. *Biological Conservation*, 2009.
- Tabibian, B., Upadhyay, U., De, A., Zarezade, A., Schölkopf, B., and Gomez-Rodriguez, M. Enhancing human learning via spaced repetition optimization. *Proceedings of the National Academy of Sciences*, 116(10): 3988–3993, 2019.
- Whitehill, J. and Movellan, J. Approximately optimal teaching of approximately optimal learners. *IEEE Transactions on Learning Technologies*, 11(2):152–164, 2017.
- Wu, L., Tian, F., Xia, Y., Fan, Y., Qin, T., Jian-Huang, L., and Liu, T.-Y. Learning to teach with dynamic loss functions. *Advances in Neural Information Processing Systems*, 31, 2018.
- Yudelson, M. V., Koedinger, K. R., and Gordon, G. J. Individualized bayesian knowledge tracing models. In *International conference on artificial intelligence in education*, pp. 171–180. Springer, 2013.
- Zhou, D., He, J., and Gu, Q. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pp. 12793–12802. PMLR, 2021.
- Zhu, X. Machine teaching: An inverse problem to machine learning and an approach toward optimal education. In *AAAI*, pp. 4083–4087, 2015.
- Zhu, X. An optimal control view of adversarial machine learning. *CoRR*, abs/1811.04422, 2018.
- Zilles, S., Lange, S., Holte, R., and Zinkevich, M. Models of cooperative teaching and learning. *JMLR*, 12(Feb): 349–384, 2011.

A. Appendix

In the appendix, we present the proofs of our theorems. The proofs of Lemma 1, Theorem 1, Theorem 2 and Theorem 3 can be found in sections C, D, E and F, respectively. In section G, we provide a reference to the existing lemmas that we rely on.

B. Extended Backgrounds on Various Learner's Models

Version-space learner. The version space learner was studied by Goldman & Kearns (1995) for machine teaching. The hypothesis class of the version space learner is usually a finite set \mathcal{H} , which contains a target hypothesis $h^* \in \mathcal{H}$. The teacher can pick a teaching example from the ground set \mathcal{X} to teach the learner. Once an example $(x, h^*(x))$ is provided to the learner, the learner will update its version space by removing those hypotheses that are not consistent with the example, i.e., $\mathcal{H} \leftarrow \mathcal{H} \setminus \{h \in \mathcal{H} | h(x) \neq h^*(x)\}$. Under this teacher-learner interaction protocol, the teacher *knows* the aforementioned update rule of the learner. The entire problem is essentially a set cover problem, which is NP-hard. But a greedy-approximation algorithm admits a teaching complexity of $\mathcal{O}(\log(\mathcal{H}) \cdot C_*)$, where C_* is the optimal teaching cost. The version space learner can also be regarded as a tabular case of the machine teaching problem, which falls in the Markov learner case, i.e., a special case of our teaching framework.

Black-box version-space learner. The black-box version-space learner was studied in Dasgupta et al. (2019). In this framework, they assume the teacher does not know the hypothesis class \mathcal{H} at the beginning, but the teacher knows the learner's dynamics rule (i.e., how does the learner update the knowledge state). Then the teaching problem is equivalent to the *online set cover* problem. The analysis of the online set cover applies to the analysis of the teaching cost. This work can be regarded as a complement to our work, as they assume the learner's model is known, but the state is unknown. Our work assumes the learner's state is observable, but the learner's model is unknown.

Black-box iterative learner. The black-box iterative learner Liu et al. (2017) is in the same philosophy as Dasgupta et al. (2019). The main difference is that, for the black-box iterative learner, it deals with gradient-based learner, i.e., the learner updates it by following the gradient descent rule. Therefore, this work still assumes the learner's model is known.

Memory-based learner. The memory-based learner was studied in Hunziker et al. (2019) for modeling the forgetting behavior of human learning. In their work, they used the half-life model as a proxy to model the human learner's model. The teaching problem can be reduced to a submodular maximization problem (maximizing the memorization utility of the underlying learner) due to the property of the half-life model.

Bayesian knowledge tracing (BKT) learner. As an instance of skill-based learners (Whitehill & Movellan, 2017), BKT assumes that student knowledge is represented as a set of binary variables, one per skill, where the skill is either mastered by the student or not. Observations in BKT are also binary: a student gets a problem/step either right or wrong. The learner's state is updated by Bayes rule given the new observation. Hence, the teacher still knows the learner's model.

C. Proof of Lemma 1

Lemma 1 Under Assumptions 1 and 2, for any $t \ge 1$, with probability at least $1 - \delta$, we have that the true parameter θ^{ℓ} lies in

$$C_t = \left\{ \boldsymbol{\theta} \in \mathbb{R}^d \,\middle| \, \|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}\|_{\boldsymbol{\Sigma}_t} \le r_t \right\},\tag{10}$$

where
$$r_t := C\sqrt{d\log((4(t^2 + t^3C^2/\lambda))/\delta)} + \sqrt{\lambda}\delta_0$$
.

Proof: We prove this by induction on k, which is the index of the value functions returned by EVI. By definition, the fitted value function in the interval $[t_k, t_{k+1} - 1]$ is $V_k(\cdot)$. To be noted, since when t = 0, we must have $\theta^{\ell} \in \mathcal{C}_0$ by Assumption 2. Therefore, we abuse the notation a little bit by reloading $t_0 = 1$ for the proof. Therefore, we first prove the base step, where $t \in [1, t_1 - 1]$. For notation simplicity, we define

$$oldsymbol{\phi}_m = oldsymbol{\phi}_{V_0}(oldsymbol{h}_m, oldsymbol{x}_m), \; oldsymbol{\Phi}_t = (oldsymbol{\phi}_1, ..., oldsymbol{\phi}_t), \; oldsymbol{v}_t = (V_0(oldsymbol{h}_2), ..., V_0(oldsymbol{h}_{t+1}))^{ op}$$
 .

Recall the definition of $\hat{\theta}_t$, by rewriting it in the matrix form, we get

$$egin{aligned} \hat{oldsymbol{ heta}}_t &= oldsymbol{\Sigma}_t^{-1} oldsymbol{b}_t = oldsymbol{\Sigma}_t^{-1} oldsymbol{b}_t + oldsymbol{\Sigma}_t^{-1} oldsymbol{b}_t + oldsymbol{\Phi}_t oldsymbol{\Phi}_t oldsymbol{b}_t \\ &= ilde{\left(\lambda oldsymbol{I} + oldsymbol{\Phi}_t oldsymbol{\Phi}_t^{ op} ig)^{-1} oldsymbol{\Phi}_t oldsymbol{v}_t - oldsymbol{\Phi}_t^{ op} oldsymbol{ heta}_0) + oldsymbol{ heta}_0. \\ &= oldsymbol{\Sigma}_t^{-1} oldsymbol{\Phi}_t oldsymbol{v}_t - oldsymbol{\Phi}_t^{ op} oldsymbol{ heta}_0) + oldsymbol{ heta}_0. \end{aligned}$$

Next, we define the following random variables

$$\eta_m = V_0(s_{m+1}) - \langle \boldsymbol{\phi}_m, \boldsymbol{\theta}^{\ell} \rangle, \quad \boldsymbol{\eta}_t = (\eta_1, ..., \eta_t)^{\top}.$$

Since $C \geq C_{\star}$, the sequence $\{\eta_t\}_{t=1}^{t_1}$ are C-sub-Gaussian. Now, we can rewrite $\hat{\theta}_t$ as

$$egin{aligned} \hat{oldsymbol{ heta}}_t &= oldsymbol{\Sigma}_t^{-1} oldsymbol{\Phi}_t \left(oldsymbol{\eta}_t + oldsymbol{\Phi}_t^ op (oldsymbol{ heta}^t - oldsymbol{ heta}_0)
ight) + oldsymbol{ heta}_0 \ &= oldsymbol{\Sigma}_t^{-1} oldsymbol{\Phi}_t oldsymbol{\eta}_t + oldsymbol{\Sigma}_t^{-1} oldsymbol{\Phi}_t oldsymbol{\Phi}_t^ op (oldsymbol{ heta}^\ell - oldsymbol{ heta}_0) + oldsymbol{ heta}_0. \end{aligned}$$

By subtracting θ^{ℓ} on both sides, we get

$$\begin{split} \hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^\ell &= \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Phi}_t \boldsymbol{\eta}_t + \left(\boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Phi}_t \boldsymbol{\Phi}_t^\top - \boldsymbol{I} \right) (\boldsymbol{\theta}^\ell - \boldsymbol{\theta}_0) \\ &= \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Phi}_t \boldsymbol{\eta}_t + \boldsymbol{\Sigma}_t^{-1} \left(\boldsymbol{\Phi}_t \boldsymbol{\Phi}_t^\top - \boldsymbol{\Sigma}_t \right) (\boldsymbol{\theta}^\ell - \boldsymbol{\theta}_0) \\ &= \boldsymbol{\Sigma}_t^{-1} \boldsymbol{\Phi}_t \boldsymbol{\eta}_t + \lambda \boldsymbol{\Sigma}_t^{-1} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\ell). \end{split}$$

Then, we further obtain the following by the Cauchy-Schwarz inequality,

$$\begin{split} \left\| \hat{\theta}_{t} - \boldsymbol{\theta}^{\ell} \right\|_{\boldsymbol{\Sigma}_{t}}^{2} &= \left\langle \boldsymbol{\Sigma}_{t} (\hat{\theta}_{t} - \boldsymbol{\theta}^{\ell}), \boldsymbol{\Phi}_{t} \boldsymbol{\eta}_{t} \right\rangle_{\boldsymbol{\Sigma}_{t}^{-1}} + \lambda \left\langle \boldsymbol{\Sigma}_{t} (\hat{\theta}_{t} - \boldsymbol{\theta}^{\ell}), \boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\ell} \right\rangle_{\boldsymbol{\Sigma}_{t}^{-1}} \\ &\leq \left\| \boldsymbol{\Sigma}_{t} (\hat{\theta}_{t} - \boldsymbol{\theta}^{\ell}) \right\|_{\boldsymbol{\Sigma}_{t}^{-1}} \left(\left\| \boldsymbol{\Phi}_{t} \boldsymbol{\eta}_{t} \right\|_{\boldsymbol{\Sigma}_{t}^{-1}} + \lambda \left\| \boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\ell} \right\|_{\boldsymbol{\Sigma}_{t}^{-1}} \right) \\ &= \left\| \hat{\boldsymbol{\theta}}_{t} - \boldsymbol{\theta}^{\ell} \right\|_{\boldsymbol{\Sigma}_{t}} \left(\left\| \boldsymbol{\Phi}_{t} \boldsymbol{\eta}_{t} \right\|_{\boldsymbol{\Sigma}_{t}^{-1}} + \lambda \left\| \boldsymbol{\theta}_{0} - \boldsymbol{\theta}^{\ell} \right\|_{\boldsymbol{\Sigma}_{t}^{-1}} \right). \end{split}$$

By Lemma 6 from Abbasi-Yadkori et al. (2011), for any $t \in [1, t_1]$, we have the following hold with probability at least $1 - \delta/(t_1(t_1 + 1))$,

$$\begin{split} \| \boldsymbol{\Phi}_t \boldsymbol{\eta}_t \|_{\boldsymbol{\Sigma}_t^{-1}} &\leq C \sqrt{2 \log \left(\frac{\det(\boldsymbol{\Sigma}_t)^{1/2}}{\lambda^{d/2} \cdot \delta / (t_1(t_1+1))} \right)} \\ &\leq C \sqrt{2 \log \left(\frac{(\lambda + tC^2)^{d/2}}{\lambda^{d/2} \cdot \delta / (t_1(t_1+1))} \right)} \\ &\leq C \sqrt{d \log \left(\frac{1 + tC^2 / \lambda}{\delta / (t_1(t_1+1))} \right)} \\ &= C \sqrt{d \log \left(\frac{t_1(t_1+1) + t \cdot t_1(1+t_1)C^2 / \lambda}{\delta} \right)}. \end{split}$$

In the next, we bound $\|oldsymbol{ heta}_0 - oldsymbol{ heta}^\ell\|_{oldsymbol{\Sigma}_t^{-1}},$

$$\left\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\ell}\right\|_{\boldsymbol{\Sigma}_t^{-1}}^2 \leq \frac{1}{\lambda_{\min}(\boldsymbol{\Sigma}_t)} \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\ell}\|_2^2 = \frac{1}{\lambda} \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\ell}\|_2^2.$$

Finally, by plugging in the above bounds, we get the desired result for the base step

$$\left\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^{\ell}\right\|_{\boldsymbol{\Sigma}_t} \leq C\sqrt{d\log\left(\frac{t_1(t_1+1) + t \cdot t_1(1+t_1)C^2/\lambda}{\delta}\right)} + \sqrt{\lambda}\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^{\ell}\|_2.$$

Since $2t \ge t_1$, then we have

$$\left\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^{\ell}\right\|_{\boldsymbol{\Sigma}_t} \le C\sqrt{d\log\left(\frac{4(t^2 + t^3C^2)/\lambda}{\delta}\right)} + \sqrt{\lambda}\delta_0.$$
 (15)

Let's suppose that, for any $k \in \{0, ..., n-1\}$, equation 15 holds for all $t \in [t_k, t_{k+1} - 1]$. For the induction step, we define the following notations for any $k \in \{0, ..., n-1\}$,

$$\tilde{V}_k(\cdot) = \min \{C, V_k(\cdot)\}.$$

Consequently, for any $k \in \{0,...,n\}$ and $t \in [t_k, t_{k+1} - 1]$, we further define

$$\breve{\boldsymbol{\Sigma}}_t = \lambda \boldsymbol{I} + \sum_{i=1}^t \boldsymbol{\phi}_{\breve{V}_{k(i)}}(\boldsymbol{h}_i, \boldsymbol{x}_i) \boldsymbol{\phi}_{\breve{V}_{k(i)}}(\boldsymbol{h}_i, \boldsymbol{x}_i)^\top, \ \breve{\boldsymbol{\mu}}_t = \lambda \boldsymbol{\theta}_0 + \sum_{i=1}^t \boldsymbol{\phi}_{\breve{V}_{k(i)}}(\boldsymbol{h}_i, \boldsymbol{x}_i) \breve{V}_{k(i)}(\boldsymbol{h}_{i+1}),$$

In analogy, we reload the definition for $\hat{\theta}_t$ and η_t by

$$reve{m{ heta}_t} = reve{m{\Sigma}}_t^{-1}reve{m{\mu}}_t, \; \eta_t = reve{V}_{k(t)}(m{h}_{t+1}) - \langle m{\phi}_{reve{V}_{k(t)}}(m{h}_t,m{x}_t), m{ heta}^\ell
angle$$

By the above definition, it's easy to verify that $\{\check{\eta}_t\}_{t=1}^{t_n}$ is almost surely C-sub-Gaussian.³ Then, we can apply the Lemma 6 again, and conclude that $\boldsymbol{\theta}^\ell \in \check{\mathcal{C}}_t$ holds with probability at least $1 - \delta/(t_n(t_n+1))$ for any $t \in [t_n, t_{n+1}-1]$ with

$$\breve{C}_t = \left\{ \boldsymbol{\theta} \in \mathbb{R}^d \left| \| \breve{\boldsymbol{\theta}}_t - \boldsymbol{\theta} \|_{\breve{\boldsymbol{\Sigma}}_t} \leq C \sqrt{d \log \left(\frac{4(t^2 + t^3 C^2 / \lambda)}{\delta} \right)} + \sqrt{\lambda} \delta_0 \right\}.$$

By the optimism principle in Algorithm 2 and the base step of induction, we will have $\check{V}_k(\cdot) = \check{V}_k$ for $k \in \{0, ..., n-1\}$, which further gives us that $\check{\Sigma}_t = \Sigma_t$, $\check{\mu}_t = \mu_t$, $\check{\eta}_t = \eta_t$ and $\check{\theta}_t = \hat{\theta}_t$ for all $t \in [1, t_{n+1} - 1]$. Consequently, we further have $\check{C}_t = \mathcal{C}_t$. Lastly, by applying the union bound over $k \geq 0$, we will get that the probability of the event in Lemma 1 holds is at least

$$1 - \sum_{k=0}^{\infty} \frac{\delta}{t_k(t_k + 1)} \ge 1 - \delta.$$

D. Proof of Theorem 1

Theorem 1 Under Assumptions 1 and 2, if the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(C_\star)$, $\lambda = 1/\delta_0^2$, and the cost function is bounded from below by c_{\min} for all non-goal knowledge states $(\mathcal{H} \setminus \{h^\star\})$ and teaching instruction (\mathcal{X}) pairs, then with probability at least $1 - 2\delta$, the teaching cost of Algorithm 1 for non-epiphany learners (i.e., $\gamma = 1$) is upper bounded by

$$\mathcal{O}\left(\left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}d\delta_{0}}{\delta c_{\min}}\right)}\right) \cdot \log^{1.5}\left(\frac{C_{\star}d}{c_{\min}\delta}\right) \cdot \frac{C_{\star}^{2}d}{c_{\min}}\right). \tag{11}$$

Proof: To prove Theorem 1, we first bound the teaching cost for running Algorithm 1 for T steps. Then, we can derive a bound for T, and plugging it back to obtain the final result.

For any T, we can decompose the teaching cost into the following

$$\sum_{t=0}^{T} c(\mathbf{h}_t, \mathbf{x}_t) \le \sum_{t=0}^{T} c(\mathbf{h}_t, \mathbf{x}_t) - V_0(\mathbf{h}_0) + C.$$
(16)

³To be noted, without such construction, if the induction step conditions on the base step, there is no guarantee that the (conditional) distribution of η_t is C-sub-Gaussian. This may prevent us from applying the Lemma 6.

By Lemma 3, we know that

$$-\sum_{t=1}^{T} \left(V_{k(t)}(\boldsymbol{h}_{t}) - V_{k(t)}(\boldsymbol{h}_{t+1}) \right) + 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + V_{0}(\boldsymbol{h}_{0}) \ge 0.$$

By adding it to the r.h.s of equation 16, we get

$$\begin{split} \sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) &\leq \sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - \underline{Y_{\boldsymbol{\theta}}}(\boldsymbol{h}_{\boldsymbol{\theta}}) + \sum_{t=1}^{T} \left(V_{k(t)}(\boldsymbol{h}_{t+1}) - V_{k(t)}(\boldsymbol{h}_{t}) \right) \\ &+ 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + \underline{Y_{\boldsymbol{\theta}}}(\boldsymbol{h}_{\boldsymbol{\theta}}) + C. \end{split}$$

By rearranging the above terms, we can get the following terms

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq \underbrace{\sum_{t=0}^{T} \left[c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - V_{k(t)}(\boldsymbol{h}_{t}) \right]}_{(1)} + \underbrace{\sum_{t=0}^{T} \left[V_{k(t)}(\boldsymbol{h}_{t+1}) - \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \right]}_{(2)} + 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + C.$$

In the next, it remains to bound (1) and (2). By Lemma 4, we can bound (1) by

Then, by Lemma 9, we can bound the martingale difference (2), with probability at least $1 - \delta$, by

$$2 \le 2C\sqrt{2T\log\left(\frac{T}{\delta}\right)}.$$

By merging the terms, we simplify the upper bound of the teaching cost to be

$$\sum_{t=1}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq 4\beta_{T} \sqrt{2Td \cdot \log\left(1 + \frac{TC^{2}}{\lambda}\right)} + 15Cd\log\left(1 + \frac{TC^{2}}{\lambda}\right) + 2C\sqrt{2T\log\left(\frac{T}{\delta}\right)},$$

where $\beta_T = C\sqrt{d\log\left((4(T^2+T^3C^2/\lambda))/\delta\right)} + \sqrt{\lambda}\delta_0$. Now, it remains to bound T. Since the cost function is bounded from below by c_{\min} , then we will have

$$T \cdot c_{\min} \leq \sum_{t=0}^{T} c(\boldsymbol{h}_t, \boldsymbol{x}_t).$$

By replacing the r.h.s. term with the upper bound derived, we get

$$T \cdot c_{\min} - C \le 4\beta_T \sqrt{2Td \cdot \log\left(1 + \frac{TC^2}{\lambda}\right)} + 15Cd\log\left(1 + \frac{TC^2}{\lambda}\right) + 2C\sqrt{2T\log\left(\frac{T}{\delta}\right)}.$$

For the terms on the r.h.s, we can loosely bound the second term by

$$15Cd\log\left(1 + \frac{TC^2}{\lambda}\right) \le 8\beta_T \sqrt{2Td \cdot \log\left(1 + \frac{TC^2}{\lambda}\right)}.$$

Then, we can bound T by

$$T \leq \frac{1}{c_{\min}} \left(12\beta_T \sqrt{2d\log\left(1 + \frac{TC^2}{\lambda}\right)} + 2C\sqrt{2\log\left(\frac{T}{\delta}\right)} \right) \cdot \sqrt{T} + \frac{C}{c_{\min}}.$$

Using the fact that $c \le a\sqrt{c} + b \Rightarrow c \le (a + \sqrt{b})^2$ for $a, b \ge 0$, we have

$$T \leq \left(\frac{1}{c_{\min}^2} \left(12\beta_T \sqrt{2d\log\left(1 + \frac{TC^2}{\lambda}\right)} + 2C\sqrt{2\log\left(\frac{T}{\delta}\right)}\right) + \sqrt{\frac{C}{c_{\min}}}\right)^2.$$

By using the inequality $(a+b)^2 \le 2a^2 + 2b^2$ twice, we get

$$T \leq \frac{32}{c_{\min}^2} \left(36\beta_T^2 d \log \left(1 + \frac{TC^2}{\lambda} \right) + C^2 \log \left(\frac{T}{\delta} \right) \right) + \frac{2C}{c_{\min}}.$$

Plugging in the following upper bound of β_T^2

$$\beta_T^2 \le 2C^2 d \log \left(\frac{4(T^2 + T^3 C^2 / \lambda)}{\delta} \right) + 2\lambda \delta_0^2 d,$$

Then, we get

$$T \leq \frac{32}{c_{\min}^2} \left(72 \left(C^2 d^2 \log \left(\frac{4T^2 + 4T^3 C^2/\lambda}{\delta}\right) + \lambda \delta_0^2 d^2\right) \cdot \log \left(1 + \frac{TC^2}{\lambda}\right) + C^2 \log \left(\frac{T}{\delta}\right)\right) + \frac{2C}{c_{\min}}.$$

By rearranging the terms, we get

$$\begin{split} T &\leq \frac{2304C^2d^2}{c_{\min}^2} \cdot \log\left(\frac{4T^2 + 4T^3C^2/\lambda}{\delta}\right) \cdot \log\left(1 + \frac{TC^2}{\lambda}\right) \\ &+ \frac{2304\lambda d^2C^2\delta_0^2}{c_{\min}^2} \cdot \log\left(1 + \frac{TC^2}{\lambda}\right) \\ &+ \frac{32C^2}{c_{\min}^2} \cdot \log\left(\frac{T}{\delta}\right) + \frac{2C}{c_{\min}}. \end{split}$$

Since $\lambda = 1/\delta_0^2$, we can get the following bound

$$\begin{split} T &\leq \frac{4608C^2d^2}{c_{\min}^2} \cdot \log\left(\frac{4T^2 + 4T^3C^2\delta_0^2}{\delta}\right) \cdot \log\left(1 + TC^2\delta_0^2\right) \\ &+ \frac{32C^2}{c_{\min}^2} \cdot \log\left(\frac{T}{\delta}\right) + \frac{2C}{c_{\min}}. \end{split}$$

We now consider the following cases: when $\delta_0 \leq 1/(TC^2)$, we will have, for some universal constant C_0 ,

$$T \le C_0 \left(\frac{C^2 d^2}{c_{\min}^2} \log^2 \left(\frac{T}{\delta} \right) \right).$$

When $\delta_0 > 1/(TC^2)$, we will have, for some universal constant C_1 ,

$$T \le C_1 \left(\frac{C^2 d^2}{c_{\min}^2} \log^2 \left(\frac{TC}{\delta} \right) \right).$$

According to Lemma 5, we arrive at the desired bound for T

$$T = \mathcal{O}\left(\frac{C^2 d^2}{c_{\min}^2} \log^2\left(\frac{Cd}{c_{\min}\delta}\right)\right).$$

Because $C = \mathcal{O}(C_{\star})$ and plugging in the bound for T into the original bound, we can finally get the desired bound for the teaching cost hold with probability at least $1 - 2\delta$ by further applying union bound on the two events (i.e., Lemma 1 and bounding 2),

$$\sum_{t} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) = \mathcal{O}\left(\left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}d\delta_{0}}{\delta c_{\min}}\right)}\right) \cdot \log^{1.5}\left(\frac{C_{\star}d}{c_{\min}\delta}\right) \cdot \frac{C_{\star}^{2}d}{c_{\min}}\right).$$

Lemma 2 Under the same assumptions as Theorem 1, if Algorithm 1 runs for T steps, then the total number of value function updates (i.e., the number of EVI calls) K is at most

$$K \le 2d \log \left(1 + \frac{TC^2}{\lambda}\right) + \log \left(1 + \frac{2T}{\lambda}\right).$$

Proof: The value function update can be triggered by either the determinant criteria (K_1) or the iteration criteria (K_2) . We bound each part separately.

Bounding K_1 : To bound K_1 , it suffices to bound the determinant of Σ_T . By Lemma 7, the fact that $\Sigma_0 = \lambda I$, and the Assumption 1, we have

$$\det(\mathbf{\Sigma}_T) \le (\lambda + TC^2)^d.$$

Therefore, we can immediately bound K_1 by

$$2^{K_1} \cdot \det(\mathbf{\Sigma}_0) = 2^{K_1} \cdot \lambda^d \le (\lambda + TC^2)^d$$

$$\Rightarrow K_1 \le 2d \log \left(1 + \frac{TC^2}{\lambda}\right).$$

Bounding K_2 : To bound K_2 , we can look at the criteria triggered by it, which immediately gives us that

$$(1+\lambda) \cdot 2^{K_2} \le T + \lambda$$

$$\Rightarrow K_2 \le \log\left(\frac{T+\lambda}{1+\lambda}\right) \le \log\left(1+\frac{T}{\lambda}\right).$$

Since $K = K_1 + K_2$, we can conclude that

$$K \le 2d \log \left(1 + \frac{TC^2}{\lambda}\right) + \log \left(1 + \frac{T}{\lambda}\right).$$

Lemma 3 Under the same assumptions as Theorem 1, for any T, the following holds,

$$\sum_{t=0}^{T} \left(V_{k(t)}(\boldsymbol{h}_{t}) - V_{k(t)}(\boldsymbol{h}_{t+1}) \right) \le 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + V_{0}(\boldsymbol{h}_{0}).$$

Proof: By Lemma 2, we can divide the T steps into K+1 segments, and within each segment, all the steps share the same value function. Let's denote the ending step of k_{th} segment as $t_{k+1}-1$, then we will have (by canceling out the intermediate terms)

$$\sum_{t=0}^{T} \left(V_{k(t)}(\boldsymbol{h}_t) - V_{k(t)}(\boldsymbol{h}_{t+1}) \right) = \sum_{k=0}^{K} V_{k}(\boldsymbol{h}_{t_k}) - V_{k}(\boldsymbol{h}_{t_{k+1}}).$$

By rearranging terms, we can further get

$$\begin{split} &\sum_{t=0}^{T} \left(V_{k(t)}(\boldsymbol{h}_{t}) - V_{k(t)}(\boldsymbol{h}_{t+1}) \right) \\ &= \sum_{k=0}^{K-1} \left(V_{k+1}(\boldsymbol{h}_{t_{k+1}}) - V_{k}(\boldsymbol{h}_{t_{k+1}}) \right) + \sum_{k=0}^{K-1} \left(V_{k}(\boldsymbol{h}_{t_{k}}) - V_{k+1}(\boldsymbol{h}_{t_{k+1}}) \right) + V_{K}(\boldsymbol{h}_{t_{K}}) - V_{K}(\boldsymbol{h}_{t_{K+1}}) \\ &= \sum_{k=0}^{K-1} \left(V_{k+1}(\boldsymbol{h}_{t_{k+1}}) - V_{k}(\boldsymbol{h}_{t_{k+1}}) \right) + V_{0}(\boldsymbol{h}_{t_{0}}) - V_{K}(\boldsymbol{h}_{t_{K}}) + V_{K}(\boldsymbol{h}_{t_{K}}) - V_{K}(\boldsymbol{h}_{t_{K+1}}) \\ &= \sum_{k=0}^{K-1} \left(V_{k+1}(\boldsymbol{h}_{t_{k+1}}) - V_{k}(\boldsymbol{h}_{t_{k+1}}) \right) + V_{0}(\boldsymbol{h}_{t_{0}}) - V_{K}(\boldsymbol{h}_{t_{K+1}}). \end{split}$$

Since the value function is non-negative, then we have

$$\sum_{t=1}^{T} (V_{k(t)}(\boldsymbol{h}_t) - V_{k(t)}(\boldsymbol{h}_{t+1})) \le K \cdot \max_{k} ||V_k||_{\infty} + V_0(\boldsymbol{h}_0).$$

By plugging in the upper bound of K from Lemma 2 and the upper bound of the value function, C, we finally arrive at

$$\sum_{t=1}^{T} \left(V_{k(t)}(\boldsymbol{h}_t) - V_{k(t)}(\boldsymbol{h}_{t+1}) \right) \le 2dC \log \left(1 + \frac{TC^2}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + V_0(\boldsymbol{h}_0).$$

Lemma 4 Under the same assumptions as Theorem 1, for any T, we can bound 1 by,

$$\widehat{I} = \sum_{t=0}^{T} \left[c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - V_{k(t)}(\boldsymbol{h}_{t}) \right] \\
\leq 4\beta_{T} \sqrt{2Td \cdot \log\left(1 + \frac{TC^{2}}{\lambda}\right)} + 2(C+1) \cdot \left(2d\log\left(1 + \frac{TC^{2}}{\lambda}\right) + \log\left(1 + \frac{T}{\lambda}\right) + 1\right),$$

where $\beta_T = C\sqrt{d\log\left((4(T^2 + T^3C^2/\lambda))/\delta\right)} + \sqrt{\lambda}\delta_0$.

Proof: First of all, by the fact that $V_{k(t)}(\boldsymbol{h}_t) = \min_{\boldsymbol{x} \in \mathcal{X}} Q_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}) = Q_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t)$, we have

Let's suppose that $Q_{k(t)}(\cdot,\cdot)$ is the value function at the $l_{k(t)}$ th value iteration of Algorithm 2, i.e., the last iteration of the while loop. Then, based on the EVI algorithm, we have

$$\begin{split} Q_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t) &= c(\boldsymbol{h}_t, \boldsymbol{x}_t) + \nu \cdot \min_{\boldsymbol{\theta} \in \mathcal{C}_t \cap \mathcal{B}} \langle \boldsymbol{\theta}, \boldsymbol{\phi}_{V^{(l_{k(t)}-1)}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle \\ &= c(\boldsymbol{h}_t, \boldsymbol{x}_t) + \nu \cdot \langle \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)}-1)}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle \\ &= c(\boldsymbol{h}_t, \boldsymbol{x}_t) + \nu \cdot \langle \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle + \nu \cdot \langle \boldsymbol{\theta}_t, [\boldsymbol{\phi}_{V^{(l_{k(t)}-1)}} - \boldsymbol{\phi}_{V^{(l_{k(t)})}}](\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle, \end{split}$$

where $\theta_t = \arg\min_{\boldsymbol{\theta} \in \mathcal{C}_t \cap \mathcal{B}} \langle \boldsymbol{\theta}, \boldsymbol{\phi}_{V^{(l_{k(t)}-1)}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle$. By plugging the above equation into \bigcirc to replace $Q_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t)$, and then rearrange terms, we get

$$\begin{split} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - Q_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \\ &= c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - \nu \cdot \langle \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \\ &- \nu \cdot \langle \boldsymbol{\theta}_{t}, [\boldsymbol{\phi}_{V^{(l_{k(t)}-1)}} - \boldsymbol{\phi}_{V^{(l_{k(t)})}}](\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \\ &= \langle \boldsymbol{\theta}^{\ell}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle - \nu \cdot \langle \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \\ &- \nu \cdot \langle \boldsymbol{\theta}_{t}, [\boldsymbol{\phi}_{V^{(l_{k(t)}-1)}} - \boldsymbol{\phi}_{V^{(l_{k(t)})}}](\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \\ &= \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle + (1 - \nu) \cdot \langle \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \\ &- \nu \cdot \langle \boldsymbol{\theta}_{t}, [\boldsymbol{\phi}_{V^{(l_{k(t)}-1)}} - \boldsymbol{\phi}_{V^{(l_{k(t)})}}](\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle. \end{split}$$

By the termination condition of the EVI algorithm, we have

$$\begin{split} &c(\boldsymbol{h}_t, \boldsymbol{x}_t) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t) - Q_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t) \\ &\leq \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle + (1 - \nu) \cdot \langle \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle + \nu \cdot \frac{1}{\lambda \cdot t'_{k(t)}} \\ &\leq \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle + (1 - \nu) \cdot C + \frac{\nu}{\lambda \cdot t'_{k(t)}}, \end{split}$$

where $t'_{k(t)}$ is the time step of $k(t)_{th}$ EVI call, we use $t'_{k(t)}$ instead of $t_{k(t)}$ to avoid ambiguity. Therefore, we can bound 1 by

$$\underbrace{1} \leq \sum_{t=0}^{T} \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle + (1 - \nu) \cdot C + \frac{\nu}{\lambda \cdot t'_{k(t)}} \\
= \sum_{t=0}^{T} \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle + \sum_{t=0}^{T} \left(\frac{\nu}{\lambda \cdot t'_{k(t)}} + (1 - \nu) \cdot C \right).$$

By the fact that both θ^{ℓ} and θ_t are in C_t and Lemma 1, we must have

$$\|\boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_t\|_{\boldsymbol{\Sigma}_t} \le 2\beta_t \le 2\beta_T.$$

Together with the Cauchy-Schwartz inequality, we obtain

$$\begin{split} \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle &\leq \|\boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}\|_{\boldsymbol{\Sigma}_{t}} \cdot \|\boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})\|_{\boldsymbol{\Sigma}_{t}^{-1}} \\ &\leq 2\|\boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}\|_{\boldsymbol{\Sigma}_{t}} \cdot \|\boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})\|_{\boldsymbol{\Sigma}_{t}^{-1}} \\ &\leq 4\beta_{T} \|\boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})\|_{\boldsymbol{\Sigma}^{-1}} \end{split}$$

In the meantime, we also have

$$\langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_t, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_t, \boldsymbol{x}_t) \rangle \leq C.$$

Then, since $C \leq \beta_T$, we get

$$\langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle \leq \min \left\{ C, 4\beta_{T} \| \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \|_{\boldsymbol{\Sigma}_{t}^{-1}} \right\}.$$

$$\leq \min \left\{ \beta_{T}, 4\beta_{T} \| \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \|_{\boldsymbol{\Sigma}_{t}^{-1}} \right\}.$$

By Lemma 8, we have

$$\sum_{t=0}^{T} \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle$$

$$\leq 4\beta_{T} \sum_{t=0}^{T} \min \left\{ 1, \| \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \|_{\boldsymbol{\Sigma}_{t}^{-1}} \right\}$$

$$\leq 4\beta_{T} \sqrt{T \cdot \left(\sum_{t=0}^{T} \min \left\{ 1, \| \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \|_{\boldsymbol{\Sigma}_{t}^{-1}} \right\} \right)}$$

$$\leq 4\beta_{T} \sqrt{T \cdot \left[2d \log \left(\frac{\operatorname{tr}(\lambda \boldsymbol{I}) + TC^{2}d}{d} \right) - \log \det(\lambda \boldsymbol{I}) \right]}$$

$$\leq 4\beta_{T} \sqrt{2Td \cdot \log \left(1 + \frac{TC^{2}}{\lambda} \right)}.$$

Next, we will bound the other part. By plugging in $1 - \nu = 1/(\lambda \cdot t_{k(t)}')$, we have

$$\sum_{t=0}^{T} \left(\frac{\nu}{\lambda \cdot t'_{k(t)}} + (1-\nu) \cdot C \right) \leq \sum_{t=0}^{T} \frac{C+1}{\lambda \cdot t'_{k(t)}} = (C+1) \sum_{t=0}^{T} \frac{1}{\lambda \cdot t'_{k(t)}}.$$

Considering the iteration triggering criteria, we get

$$t'_{k(t)+1} \le 2t'_{k(t)} + \lambda.$$

Then, we can conclude that

$$\sum_{t=0}^{T} \left(\frac{\nu}{\lambda \cdot t'_{k(t)}} + (1 - \nu) \cdot C \right) \le \sum_{k=0}^{K} (C + 1) \cdot \left(\frac{1}{\lambda} + \frac{1}{t'_{k}} \right)$$

$$\le 2(K + 1) \cdot (C + 1)$$

$$= 2(C + 1) \cdot \left(2d \log \left(1 + \frac{TC^{2}}{\lambda} \right) + \log \left(1 + \frac{T}{\lambda} \right) + 1 \right)$$

By combining the two bounds, we get

$$(1) \le 4\beta_T \sqrt{2Td \cdot \log\left(1 + \frac{TC^2}{\lambda}\right)} + 2(C+1) \cdot \left(2d\log\left(1 + \frac{TC^2}{\lambda}\right) + \log\left(1 + \frac{T}{\lambda}\right) + 1\right).$$

Lemma 5 Suppose that $T \ge 2$, $a \ge 1$ and $T \le k \log^2(aT)$ for all large enough k. Then, there exists $\eta = \eta(a)$ such that $T \le \eta \cdot k \log^2(ak)$ for all large enough k, i.e., $T = \mathcal{O}(k \log^2(ak))$.

Proof: We prove the above lemma by contrapositive. Suppose that there doesn't exist such an η . Then, we will have, for all large enough k,

$$T \ge b_k \cdot k \log^2(ak),$$

where $\{b_k\}_{k=1}^{\infty}$ is a sequence with $\lim_{k\to+\infty}b_k=+\infty$. The above inequality also implies that

$$b_k \le \frac{T}{k \log^2(ak)} \le \frac{\log^2(aT)}{\log^2(ak)}.$$

Now, let's consider the following

$$\log^{2}(aT) \le \log^{2}(ak \cdot \log^{2}(aT)) = (\log(ak) + \log\log^{2}(aT))^{2}.$$

By the inequality $(a+b)^2 \le 2a^2 + 2b^2$, we get

$$\log^2(aT) \le 2\log^2(ak) + 2\log^2(\log^2(aT)).$$

Since $aT \geq 2$, we will have

$$\log^2(\log^2(aT)) \le \frac{1}{4}\log^2(aT) \quad \Rightarrow \quad \log^2(aT) \le 2\log^2(ak) + \frac{1}{2}\log^2(aT).$$

Therefore, we can get

$$\frac{1}{2}\log^2(aT) \le 2\log^2(ak) \quad \Rightarrow \quad b_k \le \frac{\log^2(aT)}{\log^2(ak)} \le 4,$$

which leads to a contradiction with $\lim_{k\to+\infty} b_k = +\infty$. Hence, we have $T = \mathcal{O}(k \log^2(ak))$.

E. Proof of Theorem 2

Theorem 2 Under Assumptions 1 and 2, if the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(C_\star)$, $\lambda = 1/\delta_0^2$, then with probability at least $1 - 3\delta$, the total cost incurred by running Algorithm 1 for epiphany learners with $\gamma < 1$, is upper bounded by

$$\mathcal{O}\left(C_{\star} \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}^{2} \delta_{0}^{2} \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C_{\star} \log \delta}{\delta \log \gamma}\right)}\right). \tag{12}$$

Proof: The proof for the epiphany learner case mostly follows from the proof of the non-epiphany learner case, i.e., Theorem 1. In the same way, we can still decompose the cost as in Theorem 1. The only differences are in the bound of \bigcirc in and the upper bound on T.

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq \underbrace{\sum_{t=0}^{T} \left[c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - V_{k(t)}(\boldsymbol{h}_{t}) \right]}_{1} + \underbrace{\sum_{t=0}^{T} \left[V_{k(t)}(\boldsymbol{h}_{t+1}) - \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \right]}_{2} + 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + C.$$

In analogy to Lemma 4, we can get the following bound for 1,

$$\widehat{1} = \sum_{t=0}^{T} \left[c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}) - V_{k(t)}(\boldsymbol{h}_{t}) \right]
\leq \sum_{t=0}^{T} \langle \boldsymbol{\theta}^{\ell} - \boldsymbol{\theta}_{t}, \boldsymbol{\phi}_{V^{(l_{k(t)})}}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \rangle + \sum_{t=0}^{T} \left(\frac{\gamma}{\lambda \cdot t'_{k(t)}} + (1 - \gamma) \cdot C \right)
\leq 4\beta_{T} \sqrt{2Td \cdot \log\left(1 + \frac{TC^{2}}{\lambda}\right)} + \sum_{t=0}^{T} \left(\frac{\gamma}{\lambda \cdot t'_{k(t)}} + (1 - \gamma) \cdot C \right).$$

In the following, we will bound the r.h.s term in the above equation in a similar way to the proof in Lemma 4,

$$\begin{split} &\sum_{t=0}^{T} \left(\frac{\gamma}{\lambda \cdot t_{k(t)}'} + (1 - \gamma) \cdot C \right) \\ &\leq \sum_{k=0}^{K} \gamma \cdot \left(\frac{1}{\lambda} + \frac{1}{t_{k}'} \right) + (1 - \gamma) \cdot T \cdot C \\ &\leq 2\gamma \cdot \left(2d \log \left(1 + \frac{TC^{2}}{\lambda} \right) + \log \left(1 + \frac{T}{\lambda} \right) \right) + (1 - \gamma) \cdot T \cdot C. \end{split}$$

Then, by plugging in the above bounds, we get

$$(1) \le 4\beta_T \sqrt{2Td \cdot \log\left(1 + \frac{TC^2}{\lambda}\right)} + 2\gamma \cdot \left(2d\log\left(1 + \frac{TC^2}{\lambda}\right) + \log\left(1 + \frac{T}{\lambda}\right)\right) + (1 - \gamma) \cdot T \cdot C.$$

The bound for 2 in Theorem 1 still holds with probability at least $1 - \delta$. Hence, we can merge all the terms and simply them into

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq 4\beta_{T} \sqrt{2Td \cdot \log\left(1 + \frac{TC^{2}}{\lambda}\right)} + 9Cd \log\left(1 + \frac{TC^{2}}{\lambda}\right) + 2C\sqrt{2T\log\left(\frac{T}{\delta}\right)} + (1 - \gamma) \cdot T \cdot C + C$$

$$= \mathcal{O}\left(C \cdot \left(1 + d\sqrt{\log\left(1 + TC^{2}\delta_{0}^{2}\right)}\right) \cdot \sqrt{T \cdot \log\left(TC/\delta\right)} + (1 - \gamma) \cdot T \cdot C\right)$$

In the next, it's easy to show that, with probability at least $1 - \delta$, the following holds⁴

$$T = \mathcal{O}(\log(\delta)/\log(\gamma)).$$

Lastly, since $C = \mathcal{O}(C_{\star})$, and plugging in the value of T, we have the following hold with probability at least $1 - 3\delta$ by applying the union bound over the three events (i.e., Lemma 1, bounding 2) and bounding 3),

$$\sum_{t=0} c(\boldsymbol{h}_t, \boldsymbol{x}_t) = \mathcal{O}\left(C_{\star} \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C_{\star}^2 \delta_0^2 \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C_{\star} \log \delta}{\delta \log \gamma}\right)}\right).$$

F. Proof of Theorem 3

Theorem 3 For ϵ -approximate teachable epiphany learners as defined in Definition 2, if $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\|_2 \leq \delta_0$, the cost function is bounded from above by c_{max} , the confidence set C_t is constructed according to Lemma 1 with $C = \mathcal{O}(\epsilon \gamma c_{max}/(1-\gamma)^2 + C_\star)$, and if $\lambda = 1/\delta_0^2$, then with probability at least $1-3\delta$, the teaching cost incurred by running Algorithm 1 is upper bounded by

$$\mathcal{O}\left(C \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C^2 \delta_0^2 \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C \log \delta}{\delta \log \gamma}\right)} + \frac{\epsilon \log \delta}{\log \gamma}C\right). \tag{13}$$

⁴Without the loss of generality, we assume $\log(\delta)/\log(\gamma) \ge 1$.

Proof: The proof for ϵ -approximate teachable epiphany learner also follows from the proof of Theorem 1 and Theorem 2. However, to make the similar proof work, we have to bound the maximum value of the value function under \mathcal{M}_{θ^*} . To show this, by Lemma 10 and Definition 2, we have

$$||V^{\star}(\cdot|\boldsymbol{\theta}^{\star}) - V^{\star}(\cdot)||_{\infty} \le \frac{\gamma c_{\max}\epsilon}{(1-\gamma)^2},$$

where we use $V^*(\cdot|\theta^*)$ and $V^*(\cdot)$ to denote the optimal value function under the approximate MDP \mathcal{M}_{θ^*} and the true MDP \mathcal{M}_{θ^*} respectively. Therefore, we can conclude that

$$||V^{\star}(\cdot|\boldsymbol{\theta}^{\star})||_{\infty} \leq C_{\star} + \frac{\gamma c_{\max}\epsilon}{(1-\gamma)^2}.$$

Together with the optimism principle in Algorithm 2, recall that

$$\eta_t = V_{k(t)} - \langle \boldsymbol{\phi}_{V_{k(t)}}(\boldsymbol{h}_t, \boldsymbol{x}_t), \boldsymbol{\theta}^{\star} \rangle.$$

We will have η_t is $(C_\star + \frac{\gamma c \max \epsilon}{(1-\gamma)^2})$ -sub-Gaussian. Therefore, by choosing $C = C_\star + \frac{\gamma c \max \epsilon}{(1-\gamma)^2}$ as assumed, we will have the following holds with probability at least $1-\delta$ by following the same proof as in Lemma 1,

$$\theta^{\star} \in C_t \cap \mathcal{B}$$
.

Condition on the above event, the same teaching cost decomposition in Theorem 1 still holds,

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq \underbrace{\sum_{t=0}^{T} \left[c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - V_{k(t)}(\boldsymbol{h}_{t}) \right]}_{1} + \underbrace{\sum_{t=0}^{T} \left[V_{k(t)}(\boldsymbol{h}_{t+1}) - \mathbb{P}_{\boldsymbol{\theta}^{\ell}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \right]}_{2} + 2dC \log \left(1 + \frac{TC^{2}}{\lambda} \right) + C \log \left(1 + \frac{2T}{\lambda} \right) + C.$$

To bound (1), the idea is similar to Lemma 4. Due to the model misspecification, there will be one additional term in the bound,

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - Q_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})$$

$$= \sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - \mathbb{P}_{\boldsymbol{\theta}^{\star}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\star}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - Q_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})$$

$$= \sum_{t=0}^{T} \underbrace{\left(c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) + \mathbb{P}_{\boldsymbol{\theta}^{\star}} V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) - Q_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})\right)}_{\boldsymbol{x}} + \underbrace{\left[\mathbb{P} - \mathbb{P}_{\boldsymbol{\theta}^{\star}}\right] V_{k(t)}(\boldsymbol{h}_{t}, \boldsymbol{x}_{t})}_{\boldsymbol{x}}.$$

The bound of the \clubsuit term is still the same as it in Theorem 2, and the bound for the term \blacktriangledown is

$$[\mathbb{P} - \mathbb{P}_{\theta^*}]V_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t) \leq C \cdot \epsilon.$$

By putting the two bounds together we get

The bound for 2 in Theorem 1 still holds here with probability at least $1 - \delta$. Hence, we can merge all the terms and simply them into

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) \leq 4\beta_{T} \sqrt{2Td \cdot \log\left(1 + \frac{TC^{2}}{\lambda}\right)} + 9Cd\log\left(1 + \frac{TC^{2}}{\lambda}\right) + 2C\sqrt{2T\log\left(\frac{T}{\delta}\right)} + (1 - \gamma) \cdot T \cdot C + C + \epsilon \cdot T \cdot C$$

Following the proof in Theorem 2, we have the following hold with probability at least $1 - \delta$,

$$T = \mathcal{O}(\log(\delta)/\log(\gamma)).$$

By applying the union bound for the three events (i.e., Lemma 1, bounding 2) and bounding T), and plugging in the above T and $C = \mathcal{O}(C_\star + \frac{\gamma \operatorname{cmax} \epsilon}{(1-\gamma)^2})$, we can get the final bound for the teaching cost, with probability at least $1-3\delta$,

$$\sum_{t=0}^{T} c(\boldsymbol{h}_{t}, \boldsymbol{x}_{t}) = \mathcal{O}\left(C \cdot \left(1 + d\sqrt{\log\left(1 + \frac{C^{2} \delta_{0}^{2} \log \delta}{\log \gamma}\right)}\right) \cdot \sqrt{\frac{\log \delta}{\log \gamma} \log\left(\frac{C \log \delta}{\delta \log \gamma}\right)} + \frac{\epsilon \log \delta}{\log \gamma}C\right).$$

G. Additional Theorems and Lemmas

Lemma 6 (Abbasi-Yadkori et al. (2011)) Let $\{\mathcal{F}_t\}_{t=0}^{\infty}$ be a filtration. Let $\{\eta_t\}_{t=1}^{\infty}$ be a real-valued stochastic process such that η_t is \mathcal{F}_t -measurable and η_t is conditionally B-sub-Gaussian. Let $\{\phi_t\}_{t=1}^{\infty}$ be an \mathbb{R}^d -valued stochastic process such that ϕ_t is \mathcal{F}_{t-1} -measurable. Assume that Σ is a $d \times d$ positive definite matrix. For any $t \geq 0$, define

$$oldsymbol{\Sigma}_t = oldsymbol{\Sigma} + \sum_{i=1}^t oldsymbol{\phi}_i oldsymbol{\phi}_i^ op, \quad oldsymbol{s}_t = \sum_{i=1}^t \eta_i oldsymbol{\phi}_i.$$

Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t \geq 0$,

$$\|\mathbf{\Sigma}_t^{-1/2} \mathbf{s}_t\|_2 \le B \sqrt{2 \log \left(\frac{\det(\mathbf{\Sigma}_t)^{1/2}}{\delta \cdot \det(\mathbf{\Sigma})^{1/2}} \right)}.$$

Lemma 7 (Abbasi-Yadkori et al. (2011)) Suppose that $\phi_1,...,\phi_t \in \mathbb{R}^d$ and for any $1 \leq s \leq t$, we have $\|\phi_s\| \leq L$. Let $\Sigma_t = \lambda I + \sum_{s=1}^t \phi_s \phi_s^{\top}$ for some $\lambda > 0$. Then,

$$\det(\mathbf{\Sigma}_t) \leq (\lambda + tL^2/d)^d$$
.

Lemma 8 (Abbasi-Yadkori et al. (2011)) Let $\{\phi_t\}_{t=1}^{\infty}$ be in \mathbb{R}^d , and $\|\phi_t\| \leq L$ for any t. Then, for $\Sigma_t = \lambda I + \sum_{s=1}^t \phi_s \phi_s^\top$, we will have

$$\sum_{s=1}^{t} \min \left\{ 1, \|\boldsymbol{\phi}_{s}\|_{\boldsymbol{\Sigma}_{s-1}^{-1}} \right\} \leq 2 \left[d \log \left(\frac{\operatorname{tr}(\lambda \boldsymbol{I}) + tL^{2}}{d} \right) - \log \det(\lambda \boldsymbol{I}) \right].$$

Lemma 9 (Min et al. (2021)) For a transition function \mathbb{P} , a sequence of bounded and non-negative value functions $\{V_k\}_{k=1}^K$ under \mathbb{P} , and a state action sequence $\{(\boldsymbol{h}_t, \boldsymbol{x}_t)\}_{t=1}^T$, where $\|V_k\|_{\infty} \leq C$ and $\boldsymbol{h}_{t+1} \sim \mathbb{P}[\cdot | \boldsymbol{h}_t, \boldsymbol{x}_t]$, we have the following hold with probability at least $1 - \delta$,

$$\sum_{t=0}^{T} \left[V_{k(t)}(\boldsymbol{h}_t) - \mathbb{P}V_{k(t)}(\boldsymbol{h}_t, \boldsymbol{x}_t) \right] \leq 2C \sqrt{2T \log \left(\frac{T}{\delta}\right)}.$$

Lemma 10 (Csáji & Monostori (2008)) For two discounted MDPs with discounting factor γ , if they differ only in the transition functions, denoted by \mathbb{P}_1 and \mathbb{P}_2 . If their corresponding optimal value functions are V_1^* and V_2^* , respectively, and the cost function is bounded from above by c_{max} , then

$$\|V_1^{\star} - V_2^{\star}\|_{\infty} \le \frac{\gamma c_{max}}{(1 - \gamma)^2} \|\mathbb{P}_1 - \mathbb{P}_2\|_{\infty}.$$