

# UTILIZING EXPLAINABLE REINFORCEMENT LEARNING TO IMPROVE REINFORCEMENT LEARNING: A THEORETICAL AND SYSTEMATIC FRAMEWORK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement learning (RL) faces two challenges: (1) The RL agent lacks explainability. (2) The trained RL agent is, in many cases, non-optimal and even far from optimal. To address the first challenge, explainable reinforcement learning (XRL) is proposed to explain the decision-making of the RL agent. In this paper, we demonstrate that XRL can also be used to address the second challenge, i.e., improve RL performance. Our method has two parts. The first part provides a two-level explanation for why the RL agent is not optimal by identifying the mistakes made by the RL agent. Since this explanation includes the mistakes of the RL agent, it has the potential to help correct the mistakes and thus improve RL performance. The second part formulates a constrained bi-level optimization problem to learn how to best utilize the two-level explanation to improve RL performance. In specific, the upper level learns how to use the high-level explanation to shape the reward so that the corresponding policy can maximize the cumulative ground truth reward, and the lower level learns the corresponding policy by solving a constrained RL problem formulated using the low-level explanation. We propose a novel algorithm to solve this constrained bi-level optimization problem, and theoretically guarantee that the algorithm attains global optimality. We use MuJoCo experiments to show that our method outperforms state-of-the-art baselines.

## 1 INTRODUCTION

While reinforcement learning (RL) has been implemented in a wide range of applications, it faces two significant challenges: (1) The RL agent lacks transparency due to its black-box nature. (2) It has been widely observed in the RL community (Haarnoja et al., 2018; Henderson et al., 2018; Dulac-Arnold et al., 2019; Cheng et al., 2024) that, in many cases, the trained RL agent does not achieve maximum cumulative reward (i.e., non-optimal and even far from optimal). These two challenges motivate the need to improve the transparency and the performance of the RL agent.

To address the first challenge, explainable reinforcement learning (XRL) methods are proposed to explain the decision-making of the RL agents, including learning an interpretable policy (Bastani et al., 2018; Bewley & Lawry, 2021; Verma et al., 2018), pinpointing regions in the observations that are critical for choosing certain actions (Atrey et al., 2019; Guo et al., 2021a; Puri et al., 2019), learning the reward function that is actually maximized (Xie et al., 2022), and identifying the critical states that are influential to the cumulative reward (Guo et al., 2021b; Cheng et al., 2023; Amir & Amir, 2018). These XRL methods generate various explanations that improve the transparency of the RL agent and help people build trust in the RL agent (Vouros, 2022; Milani et al., 2023).

This paper demonstrates that XRL can also be used to address the second challenge of RL, i.e., RL improvement. Given a non-optimal RL agent, we use XRL to explain why this RL agent is not optimal by finding the mistakes made by the RL agent. Since our explanation provides insights into the RL agent’s mistakes, it has the potential to help correct the mistakes and thus improve performance. Some recent works (Guo et al., 2021b; Cheng et al., 2023; 2024) also use XRL to improve the RL performance. In specific, they propose to first identify the critical states that are most influential to the cumulative reward as an explanation, and then perturb the actions (Guo et al., 2021b) or fine-tune the policy (Cheng et al., 2023; 2024) at those critical states such that the refined

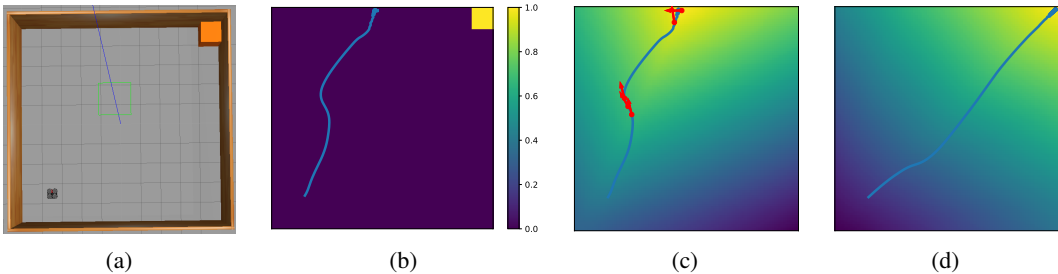
054 policy achieves higher cumulative reward. However, they do not explain why the RL agent does not  
 055 maximize the cumulative reward. This paper proposes a novel framework that first explains why the  
 056 RL agent is not optimal, and then learns how to utilize the generated explanations to improve RL  
 057 performance. We summarize our contributions as follows:

058 **Contribution statement.** This paper proposes an optimization-based framework that aims to learn  
 059 how to best utilize XRL to improve RL. We refer to this framework as “utilizing explainable RL to  
 060 improve reinforcement learning efficacy” (UTILITY). Our contributions are threefold:

061  
 062 First, we provide a two-level explanation for why the RL agent is not optimal. The high-level ex-  
 063 planation learns a reward function to which the RL agent is actually optimal, and then explains why  
 064 the RL agent is not optimal by comparing this learned reward function to the ground truth reward  
 065 function. The low-level explanation identifies the state-action pairs that lead the RL agent to be non-  
 066 optimal. We refer to these state-action pairs as “misleading” state-action pairs, and rigorously derive  
 067 a mathematical metric to identify the “misleading” state-action pairs as the low-level explanation.

068  
 069 Second, we mathematically formalize the problem of utilizing the two-level explanations to improve  
 070 the RL performance as a constrained bi-level optimization problem. In specific, the upper-level prob-  
 071 lem aims to learn how to use the high-level explanation (i.e., the learned reward function) to shape  
 072 the ground truth reward function to help the corresponding policy maximize the cumulative (ground  
 073 truth) reward. The lower-level problem learns the corresponding policy by solving a constrained RL  
 074 problem where the objective is to maximize the cumulative shaping reward and the constraint is to  
 075 discourage from visiting the low-level explanation (i.e., the “misleading” state-action pairs). Current  
 076 state-of-the-arts (Xu & Zhu, 2023; Khanduri et al., 2023) on constrained bi-level optimization can  
 077 only deal with the case where the lower-level problem is strongly convex. However, in our case, both  
 078 the objective function and constraint in the lower-level problem are highly non-convex. Therefore,  
 079 a novel theoretical framework is desired to solve this constrained bi-level optimization problem.

080  
 081 Third, we develop a novel theoretical framework and thereby an algorithm to solve the constrained  
 082 bi-level optimization problem. In specific, we first use a dual method to transform the constrained  
 083 bi-level optimization problem to an equivalent unconstrained bi-level optimization problem, and  
 084 then propose an approximation-based triple-loop algorithm to solve this unconstrained bi-level op-  
 085 timization problem. We quantify the approximation error at each loop and prove that the algorithm  
 086 attains global optimality. Experiments show that UTILITY outperforms state-of-the-art baselines.



094 Figure 1: (a) An RL task where a drone starts from the lower-left corner and navigates to the orange  
 095 goal at the upper-right corner. (b) A failing trajectory (blue) generated by the RL agent and a heat  
 096 map that visualizes the ground truth reward. (c) The two-level explanation of why the blue trajectory  
 097 fails to reach the goal. The high-level explanation is the learned reward (visualized as the heat map)  
 098 to which the RL agent is actually optimal. The low-level explanation is the misleading state-action  
 099 pairs (red circles and linked red arrows). (d) The trajectory after improvement reaches the goal.

100  
 101 **Illustrative example.** Figure 1 uses an example to illustrate our proposed framework. Suppose we  
 102 use RL to navigate a drone to the orange goal in Figure 1a. The state is the 2-D coordinate and the  
 103 action is the moving direction. The ground truth reward is one at the goal states and zero otherwise.  
 104 Figure 1b uses a heat map to visualize the ground truth reward and the blue trajectory is generated  
 105 by the learned policy. This learned policy is not optimal because it fails to reach the goal.

106  
 107 Figure 1c visualizes our two-level explanation of why the learned policy is not optimal. At the  
 high level, we use the heat map to visualize the learned reward function to which the RL agent’s

trajectory (policy) is actually optimal. The high-level explanation is that the RL agent’s policy is actually optimal to the learned reward function (visualized as the heat map in Figure 1c), and this learned reward function is very different from the ground truth reward function (visualized as the heat map in Figure 1b). Note that we normalize all the reward functions in Figures 1b-1d to  $[0, 1]$  for better comparison. For the low-level explanation, we identify the top five “misleading” state-action pairs (i.e., the red circles and arrows) in the blue trajectory where the red circles are the states and the linked red arrows are the corresponding actions chosen by the non-optimal RL agent. These state-action pairs are “misleading” since the correct actions should point to the goal.

Figure 1d shows the improvement where the heat map visualizes the learned shaping reward function and the blue trajectory is generated by the learned policy after improvement. We can see that the learned policy after improvement successfully reaches the goal.

Note that Figure 1 is just for illustration to help understand the framework. We are aware that current RL algorithms can succeed the task in Figure 1a without improvement, and we use more complicated tasks in the experiment (Section 5) to show the improvement of the proposed framework.

## 2 RELATED WORKS

Due to the space limit, we only include the related works on improving RL performance here, and we include more related works in Appendix E.

**Reward shaping.** Reward shaping can improve the RL performance by shaping the ground truth reward function. Current works on reward shaping has two main categories. The first category (Ng et al., 1999; Hu et al., 2020; Devlin & Kudenko, 2012; Gupta et al., 2022) requires an external source, such as a human expert, to provide domain knowledge as an ingredient to shape the ground truth reward function. However, when the tasks become complicated, it could be difficult and even infeasible for humans to provide domain knowledge. The second category does not need domain knowledge, including reward shaping based on exploration bonus (Bellemare et al., 2016; Ostrovski et al., 2017), learning an intrinsic reward (Zheng et al., 2018; Memarian et al., 2021), and combining exploration bonus and intrinsic reward (Devidze et al., 2022). The first category usually has better performance, while the second category does not require human-domain knowledge.

**Other methods that can improve RL performance.** Lazy-MDP (Jacq et al., 2022) shows performance improvement with the help of a provided default policy. It uses the “lazy-gap” to determine whether to choose greedy action or follow a default policy on each state  $s$ . Self-imitation learning (Oh et al., 2018) aims to encourage deep exploration by reproducing previous good decisions. Papers (Wang & Taylor, 2017; Taylor, 2018; Taylor et al., 2023) aim to improve the RL performance by utilizing external assistance, such as the assistance of a pre-trained RL agent (Wang & Taylor, 2017) or a human (Taylor, 2018; Taylor et al., 2023), which may not be accessible in some scenarios.

## 3 TWO-LEVEL EXPLANATION OF WHY THE RL AGENT IS NON-OPTIMAL

This section provides a two-level explanation to explain why the RL agent is not optimal. The RL agent’s decision making is based on a Markov decision process (MDP)  $(\mathcal{S}, \mathcal{A}, \gamma, P_0, P, r)$  which consists of a state set  $\mathcal{S}$ , an action set  $\mathcal{A}$ , a discount factor  $\gamma \in (0, 1)$ , an initial state distribution  $P_0(\cdot)$ , a state transition function  $P(\cdot|\cdot, \cdot)$ , and the ground truth reward function  $r(\cdot, \cdot)$ . The RL agent’s learned policy is denoted by  $\pi_A$  and the cumulative reward is defined as  $J_r(\pi) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$  where the initial state is drawn from  $P_0$ . When we say that the RL agent is not optimal, it means that  $\pi_A \notin \arg \max_{\pi} J_r(\pi)$ .

**The black-box assumption.** To ensure practicability, following (Bewley & Lawry, 2021; Guo et al., 2021b; Cheng et al., 2023; Guidotti et al., 2019), we only treat the RL agent as a black box with no access to its internal structure. In specific, we do not assume the access to the learned value/Q-function nor the learned policy  $\pi_A$  of the RL agent. We can only observe a set of  $m$  trajectories  $\mathcal{D} \triangleq \{\zeta^j\}_{j=1}^m$  demonstrated by the RL agent (using the non-optimal policy  $\pi_A$ ) where each trajectory  $\zeta^j = s_0, a_0, \dots$  is a state-action sequence.

**The high-level explanation.** At a high level, since the RL agent is not optimal to the ground truth reward function  $r$ , we can learn a reward function  $\hat{r}$  to which the RL agent’s policy  $\pi_A$  is actually

162 optimal, and use this learned reward function  $\hat{r}$  to generate explanations. A recent work (Xie et al.,  
 163 2022) uses the state-action pairs  $(s, a)$  with the highest  $\hat{r}(s, a)$  as an explanation, however, these  
 164 state-action pairs cannot explain why the RL agent is not optimal. Therefore, we extend (Xie et al.,  
 165 2022) by comparing the learned reward function  $\hat{r}$  to the ground truth reward function  $r$  to explain  
 166 why  $\pi_A$  is not optimal to the ground truth  $r$ . Figures 1b-1c provide an example of our high-level  
 167 explanation: the policy  $\pi_A$  is actually optimal to the learned reward function  $\hat{r}$  (visualized in Figure  
 168 1c), and this learned reward function  $\hat{r}$  is very different from the ground truth reward function  $r$   
 169 (visualized in Figure 1b).

170 Inverse reinforcement learning (IRL) (Abbeel & Ng, 2004; Ziebart et al., 2008; Arora & Doshi,  
 171 2021) can learn the reward function  $\hat{r}$  and an associated policy  $\hat{\pi}_A$  from the demonstration set  $\mathcal{D}$   
 172 such that the behaviors of policy  $\pi_A$  demonstrated in  $\mathcal{D}$  are optimal to the reward function  $\hat{r}$  learned  
 173 by IRL, and the learned policy  $\hat{\pi}_A$  can imitate the policy  $\pi_A$ . We use maximum likelihood IRL  
 174 (Zeng et al., 2022) to learn the reward function  $\hat{r}$  and policy  $\hat{\pi}_A$ .

175 While the learned reward function  $\hat{r}$  can be used for the high-level explanation, it is only interpretable  
 176 to humans in low dimension, e.g., we can use heat maps to plot reward functions (as in Figure  
 177 1). When the state and action become high dimensional, the learned reward function  $\hat{r}$  is hard for  
 178 humans to understand and thus it is difficult to straightforwardly compare  $\hat{r}$  to  $r$  (as we did in Figures  
 179 1b-1c). Therefore, we need the low-level explanation which is still interpretable in high dimension.

180 **The low-level explanation.** At a low level, the RL agent is not optimal meaning that it visits some  
 181 critical points that lead to the non-optimality. Recent works (Guo et al., 2021b; Cheng et al., 2023;  
 182 Amir & Amir, 2018; Jacq et al., 2022) identify the states that are most influential to the cumulative  
 183 reward as critical points. In order to explain why the RL agent is not optimal, we extend their idea  
 184 by redefining the critical points as the state-action pairs that lead  $\pi_A$  to be non-optimal. We refer to  
 185 these critical points as “misleading” state-action pairs and we aim to identify the top  $K$  “misleading”  
 186 state-action pairs in the demonstration set  $\mathcal{D}$  as the low-level explanation. Note that we use infinite  
 187 time-horizon MDP and it is not possible to identify the top  $K$  “misleading” state-pairs if a trajectory  
 188 has infinitely many different state-action pairs. However, in practice, the trajectory length is usually  
 189 finite and we can unify the notions of finite time horizon and infinite time horizon by introducing  
 190 “absorbing state” (Sutton & Barto, 2018). In specific, we can treat the terminal state of a finite-time-  
 191 horizon trajectory as a state keeping transitioning only to itself and generating zero reward.

192 The *key challenge* to identify the top  $K$  “misleading” state-action pairs in the demonstration set  $\mathcal{D}$  is  
 193 to propose a proper criterion or metric to define what a “misleading” state-action pair is. A straight-  
 194 forward way is to identify the state-action pairs that an optimal policy will not visit as “misleading”,  
 195 and thus use  $\pi_A(a|s) - \pi^*(a|s)$  as the metric where  $\pi^*$  is an optimal policy. However, this metric is  
 196 infeasible because the optimal policy  $\pi^*$  is not accessible.

197 In contrast, we derive a feasible metric in Definition 1 that uses a  $Q$ -function to find misleading state-  
 198 action pairs in the demonstration set  $\mathcal{D}$ . The  $Q$ -function under the policy  $\pi$  and reward function  $r$  is  
 199  $Q_r^\pi(s, a) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a]$ .

200 **Definition 1.** A state-action pair  $(s, a) \in \mathcal{D}$  is a misleading state-action pair if  $l(s, a) > 0$  where  
 201  $l(s, a) \triangleq \max_{a'} Q_r^{\pi_A}(s, a') - Q_r^{\pi_A}(s, a)$  is referred to as “misleading level”. The larger the mis-  
 202 leading level  $l(s, a)$  is, the more misleading the state-action pair  $(s, a)$  is.  
 203

204 We include the derivation of how we come up with this metric  $l$  and the proof of why  $(s, a) \in \mathcal{D}$   
 205 is misleading if  $l(s, a) > 0$  in Appendix B.1. In brief, we prove in Appendix B.1 that the policy  
 206  $\pi_A$  will be an optimal policy if  $l(s, a) = 0$  for all  $(s, a)$  such that  $a \in \pi_A(s)$  where  $\pi_A(s)$  is the  
 207 set of actions that the policy  $\pi_A$  has nonzero probability to choose at the state  $s$ . Therefore, any  
 208 state-action pair  $(s, a) \in \mathcal{D}$  such that  $l(s, a) > 0$  can be regarded as a “misleading” state-action  
 209 pair that leads the policy  $\pi_A$  to be non-optimal. The larger the misleading level  $l(s, a)$  is, the more  
 210 “misleading” the state-action pair  $(s, a)$  is, because the  $Q$  value of the chosen action  $a$  has a larger  
 211 gap from the maximum  $Q$  value at the state  $s$ . We denote the set of the identified top  $K$  “misleading”  
 212 state-action pairs by  $\mathcal{C}$ , which serves as the low-level explanation.

213 While we cannot access the policy  $\pi_A$ , we have already learned the policy  $\hat{\pi}_A$  using IRL and the  
 214 policy  $\hat{\pi}_A$  imitates the policy  $\pi_A$ . Therefore, we can use  $Q_r^{\hat{\pi}_A}$  to substitute for  $Q_r^{\pi_A}$ . Given that we  
 215 can access  $\hat{\pi}_A$  and  $r$ , we can simply learn  $Q_r^{\hat{\pi}_A}$  by sampling the environment to collect enough data  
 and doing regression. Due to the space limit, we include the method of learning  $Q_r^{\hat{\pi}_A}$  in Appendix

B.2. We are aware that in RL, it is usually sample inefficient and computationally expensive if we want to sample enough data to learn precise  $Q$ -functions corresponding to all the learning policy in the learning procedure. However, our case is different because we only need to learn one precise  $Q$ -function, which corresponds to the specific policy  $\hat{\pi}_A$ .

## 4 UTILIZING THE TWO-LEVEL EXPLANATION TO IMPROVE RL

This section provides a theoretical framework that utilizes the two-level explanation in Section 3 to improve the RL performance. In specific, Subsection 4.1 formulates the problem as a constrained bi-level optimization problem. Subsection 4.2 proposes a novel theoretical framework and thereby an algorithm to solve the constrained bi-level optimization problem.

### 4.1 PROBLEM FORMULATION

We aim to utilize the two-level explanation to improve the RL performance. For the high-level explanation  $\hat{r}$ , we use it to formulate a domain knowledge and learn how to use this domain knowledge to shape the ground truth reward  $r$  such that the learned shaping reward can lead the policy to maximize the cumulative ground truth reward  $J_r(\pi)$ . In specific, we use the comparison  $r - \hat{r}$  between the ground truth reward  $r$  and the high-level explanation  $\hat{r}$  as the domain knowledge. Note that in practice, we need to first scale  $\hat{r}$  to the same scale with  $r$  for a better comparison. Since this comparison quantifies the RL agent’s misunderstanding of the ground truth reward function  $r$ , it has the potential to help patch the error. Towards this end, we propose to learn a shaping reward function  $r_\theta(\cdot, \cdot)$  parameterized by  $\theta$ , which takes the original reward and the domain knowledge as the input. For a given state-action pair  $(s, a)$ , the corresponding shaping reward is  $r_\theta(r(s, a), r(s, a) - \hat{r}(s, a))$ .

For the low-level explanation  $\mathcal{C}$ , we discourage the RL agent from visiting the “misleading” state-action pairs in  $\mathcal{C}$ . Towards this end, we design a cost function  $c(\cdot, \cdot)$  such that  $c(s, a) \in (0, c_{\max}]$  when  $(s, a) \in \mathcal{C}$ , and  $c(s, a) = 0$  otherwise, where  $c_{\max}$  is a positive constant. We discourage from visiting  $\mathcal{C}$  by constraining the cumulative cost  $J_c(\pi) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)]$  under a budget  $b$ .

**Remark on why discouraging misleading state-action pairs can improve  $\pi_A$ .** According to monotonic policy improvement theorem,  $\pi_A$  will improve if it chooses greedy actions according to its  $Q$ -function  $Q_r^{\pi_A}$  at all the states. Given that  $\pi_A$  stops improving, it means that  $\pi_A$  must choose some nongreedy actions at some states, i.e., the misleading state-action pairs. Constraining these misleading state-action pairs means that we constrain the nongreedy actions  $\pi_A$  originally chooses at the states. This constraint can help  $\pi_A$  choose greedy actions because it eliminates some nongreedy actions and thus  $\pi_A$  only needs to find greedy actions from smaller action sets. Since this constraint can help  $\pi_A$  find greedy actions, it can help improve  $\pi_A$ .

To utilize the two-level explanation, we formulate a constrained bi-level optimization problem:

$$\max_{\theta} J_r(\pi_{r_\theta}), \text{ where } \pi_{r_\theta} = \arg \max_{\pi} \{J_{r_\theta}(\pi) + H(\pi), \text{ s.t. } J_c(\pi) \leq b\}, \quad (1)$$

where  $J_{r_\theta}(\pi) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t r_\theta(r(s_t, a_t), r(s_t, a_t) - \hat{r}(s_t, a_t))]$  is the cumulative shaping reward, and the causal entropy  $H(\pi) \triangleq E^\pi[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t)]$  is to encourage exploration and is widely used in soft  $Q$ -learning (Haarnoja et al., 2017) and soft actor-critic (Haarnoja et al., 2018).

In the problem (1), the upper level aims to learn a shaping reward function  $r_\theta$  such that the corresponding policy  $\pi_{r_\theta}$  can achieve maximum cumulative ground truth reward  $J_r(\pi_{r_\theta})$ . Given the current learned shaping reward function  $r_\theta$ , the lower-level problem in (1) is to compute the corresponding policy  $\pi_{r_\theta}$  by solving a constrained RL problem. The constrained RL problem encourages  $\pi_{r_\theta}$  to maximize the entropy-regularized cumulative shaping reward ( $J_{r_\theta}(\pi) + H(\pi)$ ) and discourages  $\pi_{r_\theta}$  from visiting  $\mathcal{C}$  by controlling the cumulative cost  $J_c(\pi)$  under the budget  $b$ . Note that if we choose  $b = 0$ , it means that the policy  $\pi_{r_\theta}$  should totally avoid the set  $\mathcal{C}$ .

Before solving the problem (1), we need to first make sure that the problem (1) is well-defined. In specific, since the lower-level problem in (1) is non-convex, it may have more than one optimal solution, i.e.,  $\pi_{r_\theta}$  is not unique. Therefore, given a reward parameter  $\theta$ , the corresponding upper-level objective function value  $J_r(\pi_{r_\theta})$  may not be unique as  $\pi_{r_\theta}$  is not unique. This will make the problem (1) ill-defined. The following theorem guarantees that the problem (1) is well-defined.

**Theorem 1.** *Given reward  $r_\theta$ , the optimal solution  $\pi_{r_\theta}$  of the lower-level problem in (1) is unique.*

## 4.2 THEORETICAL FRAMEWORK

While the current state-of-the-arts (Xu & Zhu, 2023; Khanduri et al., 2023) on constrained bi-level optimization can only deal with strongly convex lower-level problems, both the objective function and the constraint of the lower-level problem in (1) are non-convex. Therefore, a novel theoretical framework is desired to solve the problem (1). This subsection proposes a novel theoretical framework to solve the problem (1).

The proposed theoretical framework has three parts. (i) The first part transforms the original constrained bi-level optimization problem (1) to an equivalent unconstrained bi-level optimization problem. The benefit of this transformation is that the equivalent unconstrained bi-level optimization problem has an unconstrained and convex lower-level problem, which is more tractable and easier to solve. (ii) The second part proposes a novel algorithm to solve the problem (1) by solving the equivalent unconstrained bi-level optimization problem. (iii) The third part theoretically guarantees that the proposed algorithm attains global optimality.

### 4.2.1 PROBLEM TRANSFORMATION

The lower-level problem of the problem (1) is non-convex. To deal with the non-convexity issue, we introduce the dual function of the lower-level problem in (1):  $G(\lambda; \theta) \triangleq \max_{\pi} J_{r_{\theta}}(\pi) + H(\pi) - \lambda(J_c(\pi) - b)$  where  $\lambda$  is the dual variable. The dual function  $G(\lambda; \theta)$  is convex in  $\lambda$  since it is the point-wise maximum over a set of affine functions of  $\lambda$  (Boyd & Vandenberghe, 2004).

**Theorem 2.** *The optimal solution of the lower-level problem in (1) is uniquely the constrained soft policy  $\pi_{\lambda^*(\theta); \theta}$  where  $\lambda^*(\theta)$  is the unique optimal solution of the dual problem  $\min_{\lambda} G(\lambda; \theta)$ .*

We include the analytical expression of the constrained soft policy  $\pi_{\lambda^*(\theta); \theta}$  (Liu & Zhu, 2022) in Appendix C. Theorem 2 indicates that  $\pi_{\lambda^*(\theta); \theta}$  is the unique optimal solution of the lower-level problem in (1) (i.e.,  $\pi_{\lambda^*(\theta); \theta} = \pi_{r_{\theta}}$ ), and  $\lambda^*(\theta) = \arg \min_{\lambda} G(\lambda; \theta)$ . Therefore, we can replace  $\pi_{r_{\theta}}$  with  $\pi_{\lambda^*(\theta); \theta}$  and replace the lower-level problem in (1) with its dual problem, and thereby transform the constrained bi-level optimization problem (1) to the following unconstrained bi-level optimization problem:

$$\max_{\theta} J_r(\pi_{\lambda^*(\theta); \theta}), \text{ where } \lambda^*(\theta) = \arg \min_{\lambda} G(\lambda; \theta). \quad (2)$$

Compared to the original problem (1), the lower-level problem of the problem (2) is unconstrained and convex. However, there are still two challenges to solve the new problem (2).

Challenge (i): Evaluating the dual function  $G(\lambda; \theta)$  needs to obtain the constrained soft policy  $\pi_{\lambda; \theta} = \arg \max_{\pi} J_{r_{\theta}}(\pi) + H(\pi) - \lambda(J_c(\pi) - b)$ . However, current RL algorithms can only approach  $\pi_{\lambda; \theta}$  at a certain rate and only obtain the exact  $\pi_{\lambda; \theta}$  when iteration number goes to infinity. In practice, we can only run an algorithm for finite iterations and thus we cannot obtain the exact  $\pi_{\lambda; \theta}$ . This will cause errors when we evaluate the dual function  $G$ .

Challenge (ii): Even if we can obtain the exact  $\pi_{\lambda; \theta}$ , we cannot guarantee to get the exact optimal solution  $\lambda^*(\theta)$  of the lower-level problem in finite time. This makes it difficult to evaluate and solve the upper-level problem in (2) since the upper-level problem in (2) requires  $\lambda^*(\theta)$ .

### 4.2.2 THE PROPOSED ALGORITHM

This part proposes a novel algorithm that solves problem (1) by solving problem (2). The proposed algorithm is triple-loop where the inner loop approximates the constrained soft policy  $\pi_{\lambda; \theta}$  and tackles Challenge (i), the middle loop approximates the optimal solution  $\lambda^*(\theta)$  of the lower-level problem in (2) and tackles Challenge (ii), and the outer loop solves the upper-level problem in (2). We use  $n$ ,  $\bar{n}$ , and  $\tilde{n}$  to respectively denote the iteration indices of outer, middle, and inner loop.

Algorithm 1 first generates the two-level explanation (line 1) and then uses three loops to utilize the generated two-level explanation. In specific, the inner loop (lines 4-7) approximates the constrained soft policy  $\pi_{\lambda; \theta}$ . With the approximated policy  $\hat{\pi}_{\lambda; \theta}$  (line 8), the middle loop solves the lower-level problem in (2) via  $(\bar{N} - 1)$ -step gradient descent (line 9) to approximate the optimal solution  $\lambda^*(\theta)$ . With the approximated parameter  $\hat{\lambda}(\theta)$  (line 11), the outer loop solves the upper-level problem in (2) via  $(N - 1)$ -step gradient ascent (line 12). In the following, we elaborate each loop respectively.

**Algorithm 1** Utilizing explainable reinforcement learning to improve reinforcement learning**Input:** Demonstration set  $\mathcal{D}$ , initial shaping reward parameter  $\theta_0$ , dual parameter  $\lambda_0$ , and policy  $\pi_0$ **Output:** Shaping reward  $r_{\theta_N}$  and the policy after improvement  $\hat{\pi}_{\hat{\lambda}(\theta_N); \theta_N}$ 


---

```

1: Generate the two-level explanation  $(\hat{r}, \mathcal{C})$ 
2: for  $n = 0, \dots, N - 1$  do
3:   for  $\bar{n} = 0, \dots, \bar{N} - 1$  do
4:     for  $\tilde{n} = 0, \dots, \tilde{N}_{\bar{n}} - 1$  do
5:       Compute the constrained soft Q function  $Q_{\lambda_{\bar{n}}; \theta_n}^{\pi_{\bar{n}}}$ 
6:       Update the policy  $\pi_{\bar{n}+1}(a|s) \propto \exp(Q_{\lambda_{\bar{n}}; \theta_n}^{\pi_{\bar{n}}}(s, a))$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ 
7:     end for
8:     Set  $\hat{\pi}_{\lambda_{\bar{n}}; \theta_n} = \pi_{\tilde{N}_{\bar{n}}}$  and use  $\hat{\pi}_{\lambda_{\bar{n}}; \theta_n}$  to compute the approximated gradient  $g_{\lambda_{\bar{n}}; \theta_n}$ 
9:     Update  $\lambda_{\bar{n}+1} = \lambda_{\bar{n}} - \alpha_{\bar{n}} g_{\lambda_{\bar{n}}; \theta_n}$ 
10:   end for
11:   Set  $\hat{\lambda}(\theta_n) = \frac{1}{\bar{N}} \sum_{\bar{n}=0}^{\bar{N}-1} \lambda_{\bar{n}}$  and compute  $\hat{\pi}_{\hat{\lambda}(\theta_n); \theta_n}$  via  $(\bar{N} - 1)$ -step soft policy iteration
12:   Use  $\hat{\pi}_{\hat{\lambda}(\theta_n); \theta_n}$  to compute the approximated gradient  $g_{\theta_n}$  and update  $\theta_{n+1} = \theta_n + \beta_n g_{\theta_n}$ 
13: end for

```

---

**The inner loop.** Given the parameter  $(\lambda, \theta)$ , the inner loop aims to approximate the constrained soft policy  $\pi_{\lambda; \theta}$  via  $\tilde{N}_{\bar{n}}$ -step soft policy iteration (Haarnoja et al., 2017), and  $\tilde{N}_{\bar{n}} = \bar{n} + 1$ . Soft policy iteration has two steps: policy evaluation and policy improvement. Policy evaluation computes the constrained soft Q-function  $Q_{\lambda; \theta}^{\pi_{\bar{n}}}$  corresponding to the current policy  $\pi_{\bar{n}}$ , dual parameter  $\lambda$ , and reward parameter  $\theta$ . We include the expression of the constrained soft Q-function in Appendix C. Policy improvement aims to update the policy according to  $\pi_{\bar{n}+1}(a|s) \propto \exp(Q_{\lambda; \theta}^{\pi_{\bar{n}}}(s, a))$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . The output of the inner loop is the approximated policy  $\hat{\pi}_{\lambda; \theta} = \pi_{\tilde{N}_{\bar{n}}}$ . In practical implementations, we can update the policy  $\pi_{\bar{n}}$  via the policy update in soft Q-learning (Haarnoja et al., 2017) or actor update in soft actor-critic (Haarnoja et al., 2018). While soft Q-learning and soft actor-critic are designed for unconstrained RL, we show in Appendix C that we can revise them to approximate the constrained soft policy.

**The middle loop.** We aim to solve the lower-level problem in (2) via  $(\bar{N}-1)$ -step gradient descent.

**Lemma 1.** *The gradient of the dual function  $G$  is  $\nabla_{\lambda} G(\lambda; \theta) = b - J_c(\pi_{\lambda; \theta})$ .*

The gradient  $\nabla_{\lambda} G(\lambda; \theta)$  requires the exact constrained soft policy  $\pi_{\lambda; \theta}$  which is inaccessible. Therefore, we use the approximated policy  $\hat{\pi}_{\lambda; \theta}$  obtained from the inner loop to approximate the gradient  $\nabla_{\lambda} G(\lambda; \theta)$  via the gradient approximation  $g_{\lambda; \theta} = b - J_c(\hat{\pi}_{\lambda; \theta})$ , and solve the lower-level problem via  $(\bar{N}-1)$ -step gradient descent  $\lambda_{\bar{n}+1} = \lambda_{\bar{n}} - \alpha_{\bar{n}} g_{\lambda_{\bar{n}}; \theta}$ . The output is  $\hat{\lambda}(\theta) = \frac{1}{\bar{N}} \sum_{\bar{n}=0}^{\bar{N}-1} \lambda_{\bar{n}}$ .

**The outer loop.** We aim to solve the upper-level problem in (2) via  $(N - 1)$ -step gradient ascent. Towards this end, we generalize the  $Q$ /value function (Lin et al., 2020; Sutton & Barto, 2018). In specific, we define the  $Q$ -function of cost  $c$  under policy  $\pi$  as  $Q_c^{\pi}(s, a) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s, a_0 = a]$  and value function of cost as  $V_c^{\pi}(s) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s]$ . We define the  $Q$ -function of reward gradient  $\nabla_{\theta} r_{\theta}$  as  $Q_{\nabla_{\theta} r_{\theta}}^{\pi}(s, a) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r_{\theta}(s_t, a_t) | s_0 = s, a_0 = a]$  and value function of reward gradient as  $V_{\nabla_{\theta} r_{\theta}}^{\pi}(s) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r_{\theta}(s_t, a_t) | s_0 = s]$ . We define state-action visitation frequency as  $\psi^{\pi}(s, a) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t \mathbb{1}\{s_t = s, a_t = a\}]$  where  $\mathbb{1}\{\cdot\}$  is the indicator function.

**Lemma 2.** *The upper-level gradient is  $dJ_r(\pi_{\lambda^*(\theta); \theta})/d\theta = E_{(s,a) \sim \psi^{\pi_{\lambda^*(\theta); \theta}}}[ (Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) - C_{\pi_{\lambda^*(\theta); \theta}}(Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s))) Q_r^{\pi_{\lambda^*(\theta); \theta}}(s, a) ]$  where  $C_{\pi}$  is a constant vector if we fix policy  $\pi$ , and we include the expression of  $C_{\pi}$  in Appendix D.3.*

Since the gradient  $\frac{dJ_r(\pi_{\lambda^*(\theta); \theta})}{d\theta}$  requires the exact optimal solution  $\lambda^*(\theta)$  and the exact constrained soft policy  $\pi_{\lambda^*(\theta); \theta}$ , we can only use the policy  $\hat{\pi}_{\hat{\lambda}(\theta); \theta}$  to approximate  $\frac{dJ_r(\pi_{\lambda^*(\theta); \theta})}{d\theta}$  via  $g_{\theta} = E_{(s,a) \sim \psi^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}}[ (Q_{\nabla_{\theta} r_{\theta}}^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s) - C_{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(Q_c^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a) - V_c^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s))) Q_r^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a) ]$ . We then solve the upper-level problem in (2) via  $(N - 1)$ -step gradient ascent  $\theta_{n+1} = \theta_n + \beta_n g_{\theta_n}$ .

### 4.2.3 THEORETICAL ANALYSIS

This part quantifies the optimality of the policy after improvement  $\hat{\pi}_{\hat{\lambda}(\theta_N); \theta_N}$ . The main difficulty is that the inner loop and middle loop can only approximate the policy  $\pi_{\lambda; \theta}$  and the optimal solution  $\lambda^*(\theta)$ , and the approximation error may accumulate and ruin the convergence of the outer loop. In the following context, we sequentially quantify the convergence from the inner loop to the outer loop.

**Lemma 3** (convergence of the inner loop). *Given the parameter  $(\lambda, \theta)$ , the output  $\hat{\pi}_{\lambda; \theta}$  of the inner loop satisfies  $|\log \hat{\pi}_{\lambda; \theta}(a|s) - \log \pi_{\lambda; \theta}(a|s)| \leq O(\gamma^{\bar{N}_n})$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .*

Lemma 3 shows that inner loop converges linearly to the exact constrained soft policy  $\pi_{\lambda; \theta}$ .

**Assumption 1.** (i) *It holds that  $|r_\theta(\cdot, \cdot)| \leq C_1$  for any  $\theta$  where  $C_1$  is a positive constant. (ii) It holds that  $\|\nabla_\theta r_\theta(\cdot, \cdot)\| \leq C_2$  and  $\|\nabla_{\theta\theta}^2 r_\theta(\cdot, \cdot)\| \leq C_3$ , where  $C_2$  and  $C_3$  are some positive constants.*

Assumption 1 assumes that  $r_\theta$  is bounded, Lipschitz continuous, and smooth to  $\theta$ , which is a common assumption in RL (Wang et al., 2019; Kumar et al., 2023; Zhang et al., 2020). We next quantify the convergence of the middle loop:

**Lemma 4** (convergence of the middle loop). *Suppose Assumption 1 (ii) holds and let  $\alpha_{\bar{n}} = 1/(\bar{n} + 1)^{\bar{\eta}}$  where  $\bar{\eta} \in (1/2, 1)$ , the outputs  $(\hat{\lambda}(\theta), \hat{\pi}_{\hat{\lambda}(\theta); \theta})$  of the middle loop satisfy that (i)  $|\hat{\lambda}(\theta) - \lambda^*(\theta)| \leq O(1/\bar{N}^{1-\bar{\eta}})$ ; (ii)  $|\log \hat{\pi}_{\hat{\lambda}(\theta); \theta}(a|s) - \log \pi_{\lambda^*(\theta); \theta}(a|s)| \leq O(1/\bar{N}^{1-\bar{\eta}} + \gamma^{\bar{N}})$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .*

Lemma 4 shows that if the iteration  $\bar{N}$  of middle loop is sufficiently large, the approximation error of  $\lambda^*(\theta)$  and  $\pi_{\lambda^*(\theta); \theta}$  can be arbitrarily small. We next quantify the convergence of the outer loop:

**Theorem 3** (convergence of the outer loop). *Suppose Assumption 1 and the conditions in Lemma 4 hold and let  $\beta_n = 1/(n+1)^\eta$  where  $\eta \in (1/2, 1)$ , then it holds  $\frac{1}{\bar{N}} \sum_{n=0}^{\bar{N}-1} \|\nabla J_r(\pi_{\lambda^*(\theta_n); \theta_n})\|^2 \leq O(1/\bar{N}^{1-\eta} + 1/\bar{N}^{2-2\bar{\eta}} + \gamma^{2\bar{N}})$ .*

Theorem 3 shows that Algorithm 1 converges to stationarity when the iteration numbers  $N$  and  $\bar{N}$  go to infinity. When the state-action space is finite and  $r_\theta$  is linear, we have the following stronger result:

**Theorem 4** (optimality of the outer loop). *Suppose the conditions in Lemma 4 hold, the state-action space is finite, and  $r_\theta$  is linear. Let the step size  $\beta_n \leq \min\{(1-\gamma)^3/8, 1/\bar{L}\}$ , then it holds that  $\lim_{N \rightarrow \infty} \lim_{\bar{N} \rightarrow \infty} J_r(\hat{\pi}_{\hat{\lambda}(\theta_N); \theta_N}) - J_r^* = 0$  where  $J_r^*$  is the maximum value of  $J_r(\pi)$ , and  $\bar{L}$  is the smoothness constant of  $J_r(\pi_{\lambda^*(\theta); \theta})$  whose expression is in Lemma 9 in Appendix.*

Theorem 4 shows that when the state-action space is finite and  $r_\theta$  is linear, Algorithm 1 can find an optimal policy asymptotically when the iteration numbers  $N$  and  $\bar{N}$  go to infinity.

## 5 EXPERIMENT

This section provides experiment results for the proposed framework. In specific, we aim to answer the question: How does the proposed framework (UTILITY) compare to other RL improvement methods in terms of improving the RL performance. Towards this end, we introduce three RL improvement methods for comparisons. (i) **Fine-tune policy on initial states and critical states (RICE)** (Cheng et al., 2024): This method fine-tunes the policy starting at the original initial states and the states that are most influential to the cumulative reward. Note that (Guo et al., 2021b; Cheng et al., 2023) also use the most influential states to improve performance and (Cheng et al., 2024) shows performance superiority over (Guo et al., 2021b; Cheng et al., 2023), thus we pick (Cheng et al., 2024) to compare. (ii) **Self imitation learning (SIL)** (Oh et al., 2018): This method reproduces previous good decisions in order to encourage deep exploration. (iii) **Learning intrinsic reward (LIR)**: This method aims to learn an intrinsic reward  $\tilde{r}$  to formulate the shaping reward  $r + \tilde{r}$  (Zheng et al., 2018). We choose these three methods to compare because they respectively belong to three different categories: XRL method (RICE), reward shaping method (LIR), and other methods that can improve RL (SIL). We use soft actor-critic (SAC) (Haarnoja et al., 2018) as the baseline RL algorithm that all the above RL improvement methods use and improve from. We aim to show the



Table 1: Experiment results.

	SAC	UTILITY	RICE	SIL	LIR
Delayed HalfCheetah	383.45 $\pm$ 45.50	715.96 $\pm$ 42.78	456.14 $\pm$ 36.32	510.34 $\pm$ 39.28	548.28 $\pm$ 47.94
Delayed Hopper	192.90 $\pm$ 27.18	317.99 $\pm$ 19.62	232.55 $\pm$ 16.96	263.46 $\pm$ 20.72	247.27 $\pm$ 31.93
Delayed Walker2d	134.91 $\pm$ 20.80	242.63 $\pm$ 14.11	177.45 $\pm$ 20.14	172.28 $\pm$ 24.57	204.72 $\pm$ 25.99
Delayed Ant	68.11 $\pm$ 12.52	105.80 $\pm$ 14.38	77.01 $\pm$ 10.89	81.05 $\pm$ 13.43	78.23 $\pm$ 13.11

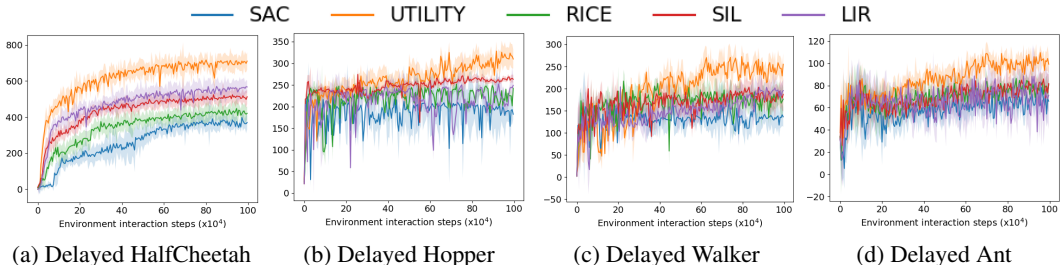


Figure 2: Improvement curve.

improvement of the above RL improvement methods compared to SAC. We also provide the results of using other baseline RL algorithms in Appendix F.1.

We test the algorithms on delayed MuJoCo environments (Zheng et al., 2018; Memarian et al., 2021; Oh et al., 2018) where the rewards are accumulated for 20 time steps and provided only at the end of these periods. Note that this makes the reward become sparse, and we include the additional experiment results on dense reward in Appendix F.2. We use four different delayed MuJoCo tasks: delayed HalfCheetah, delayed Hopper, delayed Walker2d, and delayed Ant. Following (Finn et al., 2017), each episode has the length of 100 in our experiments. We also provide the results for the case where the episode length is 1,000 in Appendix F.3.

Figure 2 shows the learning curves of the algorithms where the  $x$ -axis is the interaction steps with the MDP environment and the  $y$ -axis is the cumulative reward. We plot both the mean (i.e., the solid line) and standard deviation (i.e., the shadow area) of the algorithms. The mean and standard deviation are computed using five random seeds. From the figures, we can observe that UTILITY improves the baseline RL algorithm SAC by a large margin. While the other three methods (i.e., RICE, SIL, and LIR) can also improve SAC to some extent, UTILITY achieves the highest cumulative reward. This is due to the fact that UTILITY (i) learns a shaping reward that makes it easier to learn a good policy and (ii) discourages the policy from making the mistakes (i.e., the “misleading” state-action pairs) made by SAC. Note that the learned shaping reward is dense while the ground truth reward in the delayed MuJoCo environments is sparse, so that UTILITY can guide the learned policy to achieve higher reward. In contrast, RICE and SIL still suffer from the sparse reward. While LIR can also learn a dense shaping reward, UTILITY has the domain knowledge  $r - \hat{r}$  formulated by the high-level explanation to help better shape the reward. Moreover, UTILITY has the “misleading” state-action pairs to avoid. Note that the learned shaping reward not only helps in the sparse reward scenario, we include additional results in Appendix F.2 to show that UTILITY can still largely improve SAC when the ground truth reward is dense.

Table 1 shows the final performance of the algorithms. We can observe that UTILITY has the highest cumulative reward among all the algorithms.

**The ablation study.** Since UTILITY uses both the high-level explanation (i.e., the learned reward) to shape the ground truth reward and the low-level explanation (i.e., misleading state-action pairs) to formulate a constraint to improve SAC, we include an ablation study to separately study the effect of the shaping reward and the constraint. In specific, we consider two methods: “shaping only” and “constraint only”. The “shaping only” method only uses the high-level explanation to learn a shaping reward but does not use the low-level explanation to formulate a constraint. The “constraint only” method only uses the low-level explanation to formulate the constraint but does not shape the original reward. We include the results for the delayed environments in Table 2 and the results for the dense environments in Appendix F.4. The results in Table 2 and Appendix F.4 show that

both the “shaping only” and “constraint only” methods can improve SAC. Moreover, the “shaping only” method has a larger impact to improve the performance. This is because the shaping reward improves the policy globally as it changes the reward value for all  $(s, a)$ , while the constraint may only improve the policy locally around the misleading state-action pairs.

Table 2: Ablation study.

	SAC	UTILITY	shaping only	constraint only
Delayed HalfCheetah	383.45 ± 45.50	715.96 ± 42.78	695.63 ± 33.66	422.15 ± 22.86
Delayed Hopper	192.90 ± 27.18	317.99 ± 19.62	289.10 ± 18.41	210.12 ± 15.77
Delayed Walker2d	134.91 ± 20.80	242.63 ± 14.11	211.37 ± 18.64	175.66 ± 15.27
Delayed Ant	68.11 ± 12.52	105.80 ± 14.38	88.18 ± 8.66	75.16 ± 6.58

**Evaluation of the generated two-level explanation.** Following (Guo et al., 2021b; Cheng et al., 2023), we use fidelity as the metric to respectively evaluate the high-level and low-level explanations. The fidelity means the correctness of the two-level explanation. Since the two-level explanation is to explain why the RL agent is not optimal, one way to validate the fidelity of the explanation is to see whether the cumulative reward increases after we improve from the explanations.

From the last two columns in Table 2, we can see that both the high-level and low-level explanations are correct explanations because both the shaping only method and the constraint only method improve the performance. Moreover, the shaping only method (the fourth column in Table 2) has a higher cumulative reward than LIR (the last column in Table 1), and the constraint only method (the last column in Table 2) has a higher cumulative reward than RICE (the fourth column in Table 1). This shows the high fidelity of our two-level explanation.

To compare the fidelity of our explanation with other methods, we fix the improvement method and change the explanation to compare. For the low-level explanation, we compare with RICE. In specific, we still use the constraint only method but now the constraint is to discourage from visiting the critical states identified by RICE. We refer to this method as “RICE+constraint”. Note that it is expected that “RICE+constraint” has low fidelity in our case because RICE does not aim to explain why the RL agent is not optimal. For the high-level explanation, since there is no existing XRL method to compare, we use our shaping-only method without the domain knowledge  $r - \hat{r}$  to compare. We refer to this method as “shaping without  $r - \hat{r}$ ”. We include the results for sparse reward in Table 3 and the results for dense reward in Appendix F.5. The results show that both the high-level and low-level explanations of UTILITY have high fidelity.

Table 3: Fidelity comparison.

	SAC	shaping only (ours)	shaping without $r - \hat{r}$	constraint only (ours)	RICE+constraint
Delayed HalfCheetah	383.45 ± 45.50	695.63 ± 33.66	611.08 ± 39.44	422.15 ± 22.86	369.14 ± 19.40
Delayed Hopper	192.90 ± 27.18	289.10 ± 18.41	255.18 ± 16.57	210.12 ± 15.77	181.45 ± 12.11
Delayed Walker2d	134.91 ± 20.80	211.37 ± 18.64	191.15 ± 11.26	175.66 ± 15.27	140.26 ± 11.53
Delayed Ant	68.11 ± 12.52	88.18 ± 8.66	77.11 ± 5.52	75.16 ± 6.58	63.11 ± 4.28

**Computation time study.** While the triple-loop structure of Algorithm 1 looks computationally expensive, in practice, we can significantly accelerate the algorithm using warm start. We elaborate how we use warm start to accelerate Algorithm 1 and provide empirical results to show that Algorithm 1 is not slower than the baselines in Appendix F.6.

## 6 CONCLUSION

This paper proposes a theoretical and systematic framework that utilizes XRL to improve RL. We first provide an explanation for why the RL agent is not optimal, and then formulate the problem of utilizing the explanation to improve RL as a constrained bi-level optimization problem. We propose a novel theoretical framework to solve this problem, and use experiments to validate that the proposed framework can improve the RL performance. Despite the benefit, one limitation of the proposed algorithm is that it requires to interact with the environment. Therefore, one future work is to extend our method to the offline RL setting.

## REFERENCES

- 540  
541  
542 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In  
543 *International Conference on Machine Learning*, pp. 1–8, 2004.
- 544  
545 Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In  
546 *International Conference on Machine Learning*, pp. 22–31, 2017.
- 547  
548 Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy  
549 gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning  
550 Research*, 22(98):1–76, 2021.
- 551  
552 Dan Amir and Ofra Amir. Highlights: Summarizing agent behavior to people. In *International  
553 Conference on Autonomous Agents and MultiAgent Systems*, pp. 1168–1176, 2018.
- 554  
555 Saurabh Arora and Prashant Doshi. A survey of inverse reinforcement learning: Challenges, meth-  
556 ods and progress. *Artificial Intelligence*, 297:103500, 2021.
- 557  
558 Akanksha Atrey, Kaleigh Clary, and David Jensen. Exploratory not explanatory: Counterfactual  
559 analysis of saliency maps for deep reinforcement learning. In *International Conference on Learn-  
560 ing Representations*, 2019.
- 561  
562 Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy  
563 extraction. *Advances in Neural Information Processing Systems*, 31:2499–2509, 2018.
- 564  
565 Amir Beck. *First-order methods in optimization*. Society for Industrial and Applied Mathematics,  
566 2017.
- 567  
568 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.  
569 Unifying count-based exploration and intrinsic motivation. *Advances in Neural Information Pro-  
570 cessing Systems*, pp. 1479–1487, 2016.
- 571  
572 Tom Bewley and Jonathan Lawry. Tripletree: A versatile interpretable representation of black box  
573 agents and their environments. In *AAAI Conference on Artificial Intelligence*, volume 35, pp.  
574 11415–11422, 2021.
- 575  
576 Stephen P Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge university press, 2004.
- 577  
578 Zelei Cheng, Xian Wu, Jiahao Yu, Wenhai Sun, Wenbo Guo, and Xinyu Xing. StateMask: Ex-  
579 plaining deep reinforcement learning through state mask. In *Advances in Neural Information  
580 Processing Systems*, 2023.
- 581  
582 Zelei Cheng, Xian Wu, Jiahao Yu, Sabrina Yang, Gang Wang, and Xinyu Xing. Rice: Breaking  
583 through the training bottlenecks of reinforcement learning with explanation. In *International  
584 Conference on Machine Learning*, 2024.
- 585  
586 Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A lyapunov-  
587 based approach to safe reinforcement learning. In *Advances in Neural Information Processing  
588 Systems*, pp. 8103–8112, 2018.
- 589  
590 Rati Devidze, Parameswaran Kamalaruban, and Adish Singla. Exploration-guided reward shaping  
591 for reinforcement learning under sparse rewards. *Advances in Neural Information Processing  
592 Systems*, pp. 5829–5842, 2022.
- 593  
594 Sam Michael Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In *Internat-  
595 ional Conference on Autonomous Agents and Multiagent Systems*, pp. 433–440, 2012.
- 596  
597 Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforcement  
598 learning. *arXiv preprint arXiv:1904.12901*, 2019.
- 599  
600 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
601 of deep networks. In *International Conference on Machine Learning*, pp. 1126–1135, 2017.

- 594 Riccardo Guidotti, Anna Monreale, Fosca Giannotti, Dino Pedreschi, Salvatore Ruggieri, and  
595 Franco Turini. Factual and counterfactual explanations for black box decision making. *IEEE*  
596 *Intelligent Systems*, 34(6):14–23, 2019.
- 597
- 598 Sihang Guo, Ruohan Zhang, Bo Liu, Yifeng Zhu, Dana Ballard, Mary Hayhoe, and Peter Stone.  
599 Machine versus human attention in deep reinforcement learning tasks. *Advances in Neural Infor-*  
600 *mation Processing Systems*, 34:25370–25385, 2021a.
- 601 Wenbo Guo, Xian Wu, Usman Khan, and Xinyu Xing. Edge: Explaining deep reinforcement  
602 learning policies. *Advances in Neural Information Processing Systems*, 34:12222–12236, 2021b.
- 603
- 604 Abhishek Gupta, Aldo Pacchiano, Yuexiang Zhai, Sham Kakade, and Sergey Levine. Unpacking re-  
605 ward shaping: Understanding the benefits of reward engineering on sample complexity. *Advances*  
606 *in Neural Information Processing Systems*, 35:15281–15295, 2022.
- 607 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with  
608 deep energy-based policies. In *International Conference on Machine Learning*, pp. 1352–1361,  
609 2017.
- 610
- 611 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
612 maximum entropy deep reinforcement learning with a stochastic actor. In *International Confer-*  
613 *ence on Machine Learning*, pp. 1861–1870, 2018.
- 614 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.  
615 Deep reinforcement learning that matters. In *AAAI Conference on Artificial Intelligence*, vol-  
616 *ume 32*, 2018.
- 617
- 618 Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and  
619 Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping. *Advances*  
620 *in Neural Information Processing Systems*, 33:15931–15941, 2020.
- 621 Sandy H Huang, Kush Bhatia, Pieter Abbeel, and Anca D Dragan. Establishing appropriate trust via  
622 critical states. In *IEEE International Conference on Intelligent Robots and Systems*, pp. 3929–  
623 3936, 2018.
- 624
- 625 Alexis Jacq, Johan Ferret, Olivier Pietquin, and Matthieu Geist. Lazy-MDPs: Towards interpretable  
626 rl by learning when to act. In *International Conference on Autonomous Agents and Multiagent*  
627 *Systems*, pp. 669–677, 2022.
- 628 Prashant Khanduri, Ioannis Tsaknakis, Yihua Zhang, Jia Liu, Sijia Liu, Jiawei Zhang, and Mingyi  
629 Hong. Linearly constrained bilevel optimization: A smoothed implicit gradient approach. In  
630 *International Conference on Machine Learning*, pp. 16291–16325, 2023.
- 631
- 632 Harshat Kumar, Alec Koppel, and Alejandro Ribeiro. On the sample complexity of actor-critic  
633 method for reinforcement learning with function approximation. *Machine Learning*, pp. 1–35,  
634 2023.
- 635
- 636 Zhengxian Lin, Kin-Ho Lam, and Alan Fern. Contrastive explanations for reinforcement learning  
637 via embedded self predictions. In *International Conference on Learning Representations*, 2020.
- 638 Shicheng Liu and Minghui Zhu. Distributed inverse constrained reinforcement learning for multi-  
639 agent systems. *Advances in Neural Information Processing Systems*, 35:33444–33456, 2022.
- 640
- 641 Shicheng Liu and Minghui Zhu. Meta inverse constrained reinforcement learning: Convergence  
642 guarantee and generalization analysis. In *The Twelfth International Conference on Learning Rep-*  
643 *resentations*, 2023.
- 644
- 645 Shicheng Liu and Minghui Zhu. Learning multi-agent behaviors from distributed and streaming  
646 demonstrations. *Advances in Neural Information Processing Systems*, 36, 2024a.
- 647
- 647 Shicheng Liu and Minghui Zhu. In-trajectory inverse reinforcement learning: Learn incrementally  
from an ongoing trajectory. *arXiv preprint arXiv:2410.15612*, 2024b.

- 648 Yongshuai Liu, Jiaxin Ding, and Xin Liu. IPO: Interior-point policy optimization under constraints.  
649 In *AAAI Conference on Artificial Intelligence*, volume 34, pp. 4940–4947, 2020.  
650
- 651 Farzan Memarian, Wonjoon Goo, Rudolf Lioutikov, Scott Niekum, and Ufuk Topcu. Self-supervised  
652 online reward shaping in sparse-reward environments. In *IEEE International Conference on In-*  
653 *telligent Robots and Systems*, pp. 2369–2375, 2021.
- 654 Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. Explainable reinforcement learn-  
655 ing: A survey and comparative review. *ACM Computing Surveys*, 2023.  
656
- 657 Andrew Y Ng, Daishi Harada, and Stuart J Russell. Policy invariance under reward transformations:  
658 Theory and application to reward shaping. In *International Conference on Machine Learning*, pp.  
659 278–287, 1999.
- 660 Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In *International*  
661 *conference on machine learning*, pp. 3878–3887, 2018.  
662
- 663 Georg Ostrovski, Marc G Bellemare, Aäron Oord, and Rémi Munos. Count-based exploration with  
664 neural density models. In *International Conference on Machine Learning*, pp. 2721–2730, 2017.
- 665 Nikaash Puri, Sukriti Verma, Piyush Gupta, Dhruv Kayastha, Shripad Deshmukh, Balaji Krishna-  
666 murthy, and Sameer Singh. Explain your move: Understanding agent actions using specific and  
667 relevant feature attribution. In *International Conference on Learning Representations*, 2019.  
668
- 669 Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning by  
670 pid lagrangian methods. In *International Conference on Machine Learning*, pp. 9133–9143, 2020.
- 671 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.  
672
- 673 Matthew E Taylor. Improving reinforcement learning with human input. In *International Joint*  
674 *Conference on Artificial Intelligence*, pp. 5724–5728, 2018.
- 675 Matthew E Taylor, Nicholas Nissen, Yuan Wang, and Neda Navidi. Improving reinforcement learn-  
676 ing with human assistance: An argument for human subject studies with hippo gym. *Neural*  
677 *Computing and Applications*, 35(32):23429–23439, 2023.  
678
- 679 Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. In  
680 *International Conference on Learning Representations*, 2018.
- 681 Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri.  
682 Programmatically interpretable reinforcement learning. In *International Conference on Machine*  
683 *Learning*, pp. 5045–5054, 2018.  
684
- 685 George A Vouros. Explainable deep reinforcement learning: State of the art and challenges. *ACM*  
686 *Computing Surveys*, 55(5):1–39, 2022.
- 687 Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global  
688 optimality and rates of convergence. In *International Conference on Learning Representations*,  
689 2019.
- 690 Zhaodong Wang and Matthew E Taylor. Improving reinforcement learning with confidence-based  
691 demonstrations. In *International Joint Conference on Artificial Intelligence*, pp. 3027–3033,  
692 2017.  
693
- 694 Yuansheng Xie, Soroush Vosoughi, and Saeed Hassanpour. Towards interpretable deep reinforce-  
695 ment learning models via inverse reinforcement learning. In *International Conference on Pattern*  
696 *Recognition*, pp. 5067–5074, 2022.
- 697 Siyuan Xu and Minghui Zhu. Efficient gradient approximation method for constrained bilevel opti-  
698 mization. In *AAAI Conference on Artificial Intelligence*, volume 37, pp. 12509–12517, 2023.  
699
- 700 Tengyu Xu, Yingbin Liang, and Guanghui Lan. CRPO: A new approach for safe reinforcement  
701 learning with convergence guarantee. In *International Conference on Machine Learning*, pp.  
11480–11491, 2021.

702 Siliang Zeng, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Maximum-likelihood inverse re-  
703 inforcement learning with finite-time guarantees. *Advances in Neural Information Processing*  
704 *Systems*, 35:10122–10135, 2022.

705  
706 Kaiqing Zhang, Alec Koppel, Hao Zhu, and Tamer Basar. Global convergence of policy gradient  
707 methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):  
708 3586–3612, 2020.

709 Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On learning intrinsic rewards for policy gradient  
710 methods. *Advances in Neural Information Processing Systems*, pp. 4649–4659, 2018.

711  
712 Zhengyuan Zhou, Michael Bloem, and Nicholas Bambos. Infinite time horizon maximum causal  
713 entropy inverse reinforcement learning. *IEEE Transactions on Automatic Control*, 63(9):2787–  
714 2802, 2017.

715 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse  
716 reinforcement learning. In *AAAI Conference on Artificial Intelligence*, pp. 1433–1438, 2008.

## 717 718 719 A APPENDIX

720  
721 This appendix has three sections. Section B provides additional content of the two-level explanation.  
722 Section C provides notions and notations that serve as the building blocks of the appendix. Section  
723 D provides the proof of all the lemmas and theorems in the paper. Section F provides experiment  
724 details.

## 725 726 B TWO-LEVEL EXPLANATION

727  
728 This section has two subsections. Subsection B.1 provides the derivation steps of how we come up  
729 with the mathematical metric  $l$  to find “misleading” state-action pairs in Section 3 and proves that  
730 a state-action pair  $(s, a) \in \mathcal{D}$  is misleading if  $l(s, a) > 0$ . Subsection B.2 provides an algorithm to  
731 learn the  $Q$ -function  $Q_r^{\pi_A}$ .

### 732 733 B.1 JUSTIFICATION OF THE MATHEMATICAL METRIC TO FIND THE “MISLEADING” 734 STATE-ACTION PAIRS

735 Since the misleading state-action pairs lead the policy  $\pi_A$  to be non-optimal, we can say that the pol-  
736 icy  $\pi_A$  will be an optimal policy if it does not visit misleading state-action pairs. Therefore, in order  
737 to identify “misleading” state-action pairs using  $Q$ -functions, we need to first build a connection  
738 between  $Q$ -functions and policy:

739  
740 **Definition 2.** (i) When we say that a  $Q$ -function  $\bar{Q}_r^\pi$  **indicates** a policy  $\pi$ , it means that  $\pi(s) =$   
741  $\arg \max_a \bar{Q}_r^\pi(s, a)$  for any  $s \in \mathcal{S}$  where  $\pi(s)$  is the set of actions that the policy  $\pi$  will choose at  
742 the state  $s$ .

743 (ii) We use  $\bar{Q}_r^\pi$  to denote the set of all the  $Q$ -functions that indicate the policy  $\pi$ , and thus we can  
744 say that  $\bar{Q}_r^\pi$  **indicates** the policy  $\pi$ . Moreover, when we say that  $\bar{Q}_r^\pi$  **indicates** an action  $a$  at a state  
745  $s$ , it means that  $a \in \arg \max_{a'} \bar{Q}_r^\pi(s, a')$  where  $\bar{Q}_r^\pi$  is an arbitrary  $Q$ -function in  $\bar{Q}_r^\pi$ .

746 (iii) When we say that  $Q_r^\pi$  **disagrees with**  $\bar{Q}_r^\pi$  **on**  $(s, a)$ , it means that  $\bar{Q}_r^\pi$  indicates the action  $a$  at  
747 the state  $s$  but  $a \notin \arg \max_{a'} Q_r^\pi(s, a')$ .

748 Definition 2 establishes a connection between  $Q$ -functions and policy. With this connection, the  
749 following theorem provides a way to define a mathematical metric to find “misleading” state-action  
750 pairs using  $Q$ -functions:

751 **Theorem 5.** The policy  $\pi_A$  will be an optimal policy if  $\pi_A$  never visits any state-action pair  $(s, a)$ ,  
752 on which  $Q_r^{\pi_A}$  disagrees with  $\bar{Q}_r^{\pi_A}$ .

753  
754 *Proof.* Suppose there is no state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , on which  $Q_r^{\pi_A}$  disagrees with  
755  $\bar{Q}_r^{\pi_A}$ . Then for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , if  $a \in \arg \max_{a'} \bar{Q}_r^{\pi_A}(s, a')$  where  $\bar{Q}_r^{\pi_A} \in \bar{Q}_r^{\pi_A}$ , it holds  
that  $a \in \arg \max_{a'} Q_r^{\pi_A}(s, a')$ . Since  $\pi_A(s) = \arg \max_a \bar{Q}_r^{\pi_A}(s, a)$  where  $\bar{Q}_r^{\pi_A} \in \bar{Q}_r^{\pi_A}$ , it

holds that  $\pi_A(s) \subseteq \arg \max_a Q_r^{\pi_A}(s, a)$  for any  $s \in \mathcal{S}$ . Recall that  $V_r^{\pi_A}$  and  $Q_r^{\pi_A}$  are respectively the value function and  $Q$ -function of the policy  $\pi_A$  under the reward function  $r$ , then it holds that  $Q_r^{\pi_A}(s, a) = r(s, a) + \gamma E_{s' \sim P(\cdot|s, a)}[V_r^{\pi_A}(s')]$  and  $V_r^{\pi_A}(s) = E_{a' \sim \pi_A(\cdot|s)}[Q_r^{\pi_A}(s, a')]$ . Since  $\pi_A(s) \subseteq \arg \max_a Q_r^{\pi_A}(s, a)$  for any  $s \in \mathcal{S}$ , then  $V_r^{\pi_A}(s) = E_{a' \sim \pi_A(\cdot|s)}[Q_r^{\pi_A}(s, a')] = \max_{a'} Q_r^{\pi_A}(s, a')$  for any  $s \in \mathcal{S}$ . Therefore, the  $Q$ -function  $Q_r^{\pi_A}$  satisfies the Bellman optimality equation  $Q_r^{\pi_A}(s, a) = r(s, a) + \gamma E_{s' \sim P(s'|s, a)}[\max_{a'} Q_r^{\pi_A}(s', a')]$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and thus  $Q_r^{\pi_A}$  is the optimal  $Q$ -function because the Bellman optimality equation is uniquely satisfied by the optimal  $Q$ -function. Since  $\pi_A(s) \subseteq \arg \max_a Q_r^{\pi_A}(s, a)$  for any  $s \in \mathcal{S}$ , the policy  $\pi_A$  should be an optimal policy.  $\square$

Since the misleading state-action pairs lead the policy  $\pi_A$  to be non-optimal, we can say that the policy  $\pi_A$  will be an optimal policy if it has zero probability to visit misleading state-action pairs. Therefore, Theorem 5 shows that the “misleading” state-action pairs can be mathematically defined as the state-action pairs, on which  $Q_r^{\pi_A}$  disagrees with  $\bar{Q}_r^{\pi_A}$ . Therefore, we can develop a mathematical metric to find the top  $K$  “misleading” state-action pairs in the demonstration set  $\mathcal{D}$  using  $Q$ -functions. Since the demonstration set  $\mathcal{D}$  is generated by the policy  $\pi_A$  and  $Q_r^{\pi_A}$  indicates the policy  $\pi_A$ ,  $\bar{Q}_r^{\pi_A}$  will indicate the action  $a$  at the state  $s$  for any  $(s, a) \in \mathcal{D}$ . In order to find the “misleading” state-actions, we need to find  $(s, a) \in \mathcal{D}$  such that  $a \notin \arg \max_{a'} Q_r^{\pi_A}(s, a')$  or in other words,  $Q_r^{\pi_A}(s, a) < \max_{a'} Q_r^{\pi_A}(s, a')$ . Therefore, we can use the function  $l(s, a) \triangleq \max_{a'} Q_r^{\pi_A}(s, a') - Q_r^{\pi_A}(s, a)$  as the metric to identify “misleading” state-action pairs in the demonstration set  $\mathcal{D}$ . The larger the loss  $l(s, a)$  is, the more “misleading” the state-action pair  $(s, a)$  is, because the  $Q$  value of the chosen action  $a$  has a larger gap from the maximum  $Q$  value at the state  $s$  under the  $Q$ -function  $Q_r^{\pi_A}$ .

## B.2 METHOD

In this subsection, we use a standard regression method to learn the  $Q$ -function  $Q_r^{\hat{\pi}_A}$ . In specific, we roll out the policy  $\hat{\pi}_A$  to generate a set  $\bar{\mathcal{D}}$  of many  $(s, a)$  samples. For each  $(s, a) \in \bar{\mathcal{D}}$ , we use  $\hat{\pi}_A$  to generate many trajectories starting from  $(s, a)$  and use these trajectories to estimate  $Q_r^{\hat{\pi}_A}(s, a)$ . Since we can generate many trajectories, we can estimate the  $Q$  value  $Q_r^{\hat{\pi}_A}(s, a)$  for each  $(s, a) \in \bar{\mathcal{D}}$  quite accurately. Then we use a neural network  $Q_\omega$  parameterized by  $\omega$  to solve the following regression problem:

$$\min_{\omega} \sum_{(s, a) \in \bar{\mathcal{D}}} \|Q_\omega(s, a) - Q_r^{\hat{\pi}_A}(s, a)\|^2.$$

## C NOTIONS AND NOTATIONS

The shaping reward function  $r_\theta(r(s, a), r(s, a) - \hat{r}(s, a))$  is a function of  $r(s, a)$  and  $r(s, a) - \hat{r}(s, a)$ , and  $r(s, a)$  and  $r(s, a) - \hat{r}(s, a)$  are both functions of  $(s, a)$ . Therefore, the shaping reward function  $r_\theta$  is also a function of  $(s, a)$ . For simple notations, we use  $r_\theta(s, a)$  to denote the shaping reward of  $(s, a)$ . Given the policy  $\pi$  and the parameters  $(\lambda, \theta)$ s, the corresponding constrained soft  $Q$ -function and constrained soft value function are:

$$Q_{\lambda; \theta}^\pi(s, a) \triangleq r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_{\lambda; \theta}^\pi(s') ds',$$

$$V_{\lambda; \theta}^\pi(s) \triangleq E^\pi \left[ \sum_{t=0}^{\infty} \gamma^t (r_\theta(s_t, a_t) - \lambda c(s_t, a_t) - \log \pi(a_t|s_t)) \middle| s_0 = s \right].$$

Moreover, it can be shown (Zeng et al., 2022; Haarnoja et al., 2017; Liu & Zhu, 2024b) that  $\exp(V_{\lambda; \theta}^\pi(s)) = \int_{a \in \mathcal{A}} \exp(Q_{\lambda; \theta}^\pi(s, a)) da$ .

The constrained soft policy (Liu & Zhu, 2022; 2023) is

$$\pi_{\lambda; \theta}(a|s) = \frac{\exp(Q_{\lambda; \theta}(s, a))}{\exp(V_{\lambda; \theta}(s))}, \quad (3)$$

$$Q_{\lambda; \theta}(s, a) = r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_{\lambda; \theta}(s') ds', \quad (4)$$

$$V_{\lambda;\theta}(s) = \log \int_{a \in \mathcal{A}} \exp(Q_{\lambda;\theta}(s, a)) da. \quad (5)$$

We can obtain the constrained soft policy via soft Q-learning (Haarnoja et al., 2017) or soft actor-critic (Haarnoja et al., 2018) by treating  $r_\theta - \lambda c$  as the new reward function. We define the cumulative cost under the policy  $\pi$  starting from  $(s, a)$  as  $Q_c^\pi(s, a) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s, a_0 = a]$  and the cumulative cost starting from  $s$  as  $V_c^\pi(s) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s]$ . We define the cumulative reward under the policy  $\pi$  starting from  $(s, a)$  as  $Q_r^\pi(s, a) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a]$  and the cumulative reward starting from  $s$  as  $V_r^\pi(s) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$ . We define the cumulative reward gradient under the policy  $\pi$  starting from  $(s, a)$  as  $Q_{\nabla_\theta r_\theta}^\pi(s, a) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta r_\theta(s_t, a_t) | s_0 = s, a_0 = a]$  and the cumulative reward starting from  $s$  as  $V_{\nabla_\theta r_\theta}^\pi(s) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta r_\theta(s_t, a_t) | s_0 = s]$ . We define the state visitation frequency under a policy  $\pi$  as  $\psi^\pi(s) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t \mathbb{1}\{s_t = s\}]$  and state-action visitation frequency as  $\psi^\pi(s, a) \triangleq E^\pi[\sum_{t=0}^{\infty} \gamma^t \mathbb{1}\{s_t = s, a_t = a\}]$  where  $\mathbb{1}\{\cdot\}$  is the indicator function.

## D PROOF

This section provides the proof of all the lemmas and theorems in the paper.

**Lemma 5.** *The gradients of the constrained soft policy are respectively  $\nabla_\lambda \log \pi_{\lambda;\theta}(a|s) = V_c^{\pi_{\lambda;\theta}}(s) - Q_c^{\pi_{\lambda;\theta}}(s, a)$  and  $\nabla_\theta \log \pi_{\lambda;\theta}(a|s) = Q_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s, a) - V_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s)$ .*

*Proof.* Recall from (3), we know that  $\nabla_\lambda \log \pi_{\lambda;\theta}(a|s) = \nabla_\lambda Q_{\lambda;\theta}(s, a) - \nabla_\lambda V_{\lambda;\theta}(s)$ . Recall from (4), we know that

$$\begin{aligned} \nabla_\lambda Q_{\lambda;\theta}(s, a) &= -c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) \nabla_\lambda V_{\lambda;\theta}(s') ds', \\ &\stackrel{(a)}{=} -c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) \frac{1}{\exp(V_{\lambda;\theta}(s'))} \int_{a' \in \mathcal{A}} \nabla_\lambda \exp(Q_{\lambda;\theta}(s', a')) da' ds', \\ &= -c(s, a) + \gamma \int_{s' \in \mathcal{S}} \int_{a' \in \mathcal{A}} P(s'|s, a) \frac{\exp(Q_{\lambda;\theta}(s', a'))}{\exp(V_{\lambda;\theta}(s'))} \nabla_\lambda Q_{\lambda;\theta}(s', a') da' ds', \\ &\stackrel{(b)}{=} -c(s, a) + \gamma \int_{s' \in \mathcal{S}} \int_{a' \in \mathcal{A}} P(s'|s, a) \pi_{\lambda;\theta}(a|s) \nabla_\lambda Q_{\lambda;\theta}(s', a') da' ds', \\ &= -c(s, a) + \gamma \int_{s' \in \mathcal{S}} \int_{a' \in \mathcal{A}} P(s'|s, a) \pi_{\lambda;\theta}(a|s) \left[ -c(s', a') + \gamma \int_{s'' \in \mathcal{S}} P(s''|s', a') \nabla_\lambda V_{\lambda;\theta}(s'') ds'' \right], \end{aligned}$$

where (a) follows (5), (b) follows (3). Keep the expansion, we can see that

$$\nabla_\lambda Q_{\lambda;\theta}(s, a) = -E^{\pi_{\lambda;\theta}}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s, a_0 = a] = -Q_c^{\pi_{\lambda;\theta}}(s, a). \quad (6)$$

Similarly, we can get that:

$$\nabla_\lambda V_{\lambda;\theta}(s) = -E^{\pi_{\lambda;\theta}}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | s_0 = s] = -V_c^{\pi_{\lambda;\theta}}(s), \quad (7)$$

$$\nabla_\theta Q_{\lambda;\theta}(s, a) = E^{\pi_{\lambda;\theta}}[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta r_\theta(s_t, a_t) | s_0 = s, a_0 = a] = Q_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s, a), \quad (8)$$

$$\nabla_\theta V_{\lambda;\theta}(s) = E^{\pi_{\lambda;\theta}}[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta r_\theta(s_t, a_t) | s_0 = s] = V_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s), \quad (9)$$

Therefore, we can compute the gradients

$$\begin{aligned} \nabla_\lambda \log \pi_{\lambda;\theta}(a|s) &= \nabla_\lambda Q_{\lambda;\theta}(s, a) - \nabla_\lambda V_{\lambda;\theta}(s) = V_c^{\pi_{\lambda;\theta}}(s) - Q_c^{\pi_{\lambda;\theta}}(s, a), \\ \nabla_\theta \log \pi_{\lambda;\theta}(a|s) &= \nabla_\theta Q_{\lambda;\theta}(s, a) - \nabla_\theta V_{\lambda;\theta}(s) = Q_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s, a) - V_{\nabla_\theta r_\theta}^{\pi_{\lambda;\theta}}(s). \end{aligned}$$

□



**Lemma 6.** *The constrained soft operator  $\mathcal{T}_{\lambda;\theta}^{\text{soft}}$ :*

$$\begin{aligned} (\mathcal{T}_{\lambda;\theta}^{\text{soft}}Q)(s, a) &\triangleq r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) \log \left[ \int_{a' \in \mathcal{A}} \exp(Q(s', a')) da' \right] ds', \\ (\mathcal{T}_{\lambda;\theta}^{\text{soft}}V)(s) &\triangleq \log \left[ \int_{a \in \mathcal{A}} \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V(s') ds' \right) da \right], \end{aligned}$$

is a contraction map with constant  $\gamma$ .

*Proof.* It has been proved that  $\mathcal{T}_{\lambda;\theta}^{\text{soft}}Q$  is a contraction map with constant  $\gamma$  (Appendix A.2 in (Haarnoja et al., 2017)) if we replace  $r$  with  $r_\theta - \lambda c$ . Here we show that  $\mathcal{T}_{\lambda;\theta}^{\text{soft}}V$  is a contraction map with constant  $\gamma$ . Define a norm of  $V$  as  $\|V_1 - V_2\| = \sup_{s \in \mathcal{S}} |V_1(s) - V_2(s)|$  and suppose  $\|V_1 - V_2\| = \epsilon$ . Then we have that

$$\begin{aligned} \mathcal{T}_{\lambda;\theta}^{\text{soft}}V_1(s) &= \log \left[ \int_{a \in \mathcal{A}} \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_1(s') ds' \right) da \right], \\ &\leq \log \left[ \int_{a \in \mathcal{A}} \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) [V_2(s') + \epsilon] ds' \right) da \right], \\ &= \log \left[ \int_{a \in \mathcal{A}} \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_2(s') ds' + \gamma \epsilon \right) da \right], \\ &= \log \left[ \int_{a \in \mathcal{A}} \exp(\gamma \epsilon) \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_2(s') ds' \right) da \right], \\ &= \mathcal{T}_{\lambda;\theta}^{\text{soft}}V_2(s) + \gamma \epsilon. \end{aligned}$$

Similarly, we can get  $\mathcal{T}_{\lambda;\theta}^{\text{soft}}V_1(s) \geq \mathcal{T}_{\lambda;\theta}^{\text{soft}}V_2(s) - \gamma \epsilon$ . Therefore,  $\|\mathcal{T}_{\lambda;\theta}^{\text{soft}}V_1 - \mathcal{T}_{\lambda;\theta}^{\text{soft}}V_2\| \leq \gamma \epsilon = \gamma \|V_1 - V_2\|$ .  $\square$

**Lemma 7.** *It holds that  $Q_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s, a) \geq \mathcal{T}_{\lambda;\theta}^{\text{soft}}(Q_{\lambda;\theta}^{\pi_{\bar{n}}})(s, a)$  and  $V_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s) \geq \mathcal{T}_{\lambda;\theta}^{\text{soft}}(V_{\lambda;\theta}^{\pi_{\bar{n}}})(s)$  for any  $(s, a)$ .*

*Proof.*

$$\begin{aligned} Q_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s, a) &\stackrel{(a)}{=} r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) E_{a' \sim \pi_{\bar{n}+1}} [Q_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s', a') - \log \pi_{\bar{n}+1}(a'|s')] ds', \\ &\stackrel{(b)}{\geq} r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) E_{a' \sim \pi_{\bar{n}+1}(\cdot|s')} [Q_{\lambda;\theta}^{\pi_{\bar{n}}}(s', a') - \log \pi_{\bar{n}+1}(a'|s')] ds', \\ &= r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) \log \left[ \int_{a' \in \mathcal{A}} \exp(Q_{\lambda;\theta}^{\pi_{\bar{n}}}(s', a')) da' \right] ds', \\ &= \mathcal{T}_{\lambda;\theta}^{\text{soft}}(Q_{\lambda;\theta}^{\pi_{\bar{n}}})(s, a), \end{aligned}$$

where (a) follow equations (2)-(3) in (Haarnoja et al., 2018) and (b) follows policy improvement theorem (Theorem 4 in (Haarnoja et al., 2017)). Similarly, we can get that

$$\begin{aligned} V_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s) &= E_{a \sim \pi_{\bar{n}+1}(\cdot|s)} [Q_{\lambda;\theta}^{\pi_{\bar{n}+1}}(s, a) - \log \pi_{\bar{n}+1}(a|s)], \\ &\geq E_{a \sim \pi_{\bar{n}+1}(\cdot|s)} [Q_{\lambda;\theta}^{\pi_{\bar{n}}}(s, a) - \log \pi_{\bar{n}+1}(a|s)], \\ &= \log \left[ \int_{a \in \mathcal{A}} \exp(Q_{\lambda;\theta}^{\pi_{\bar{n}}}(s, a)) da \right], \\ &= \log \left[ \int_{a \in \mathcal{A}} \exp \left( r_\theta(s, a) - \lambda c(s, a) + \gamma \int_{s' \in \mathcal{S}} P(s'|s, a) V_{\lambda;\theta}^{\pi_{\bar{n}}}(s') ds' \right) da \right], \\ &= \mathcal{T}_{\lambda;\theta}^{\text{soft}}(V_{\lambda;\theta}^{\pi_{\bar{n}}})(s). \end{aligned}$$

$\square$

918 D.1 PROOF OF LEMMA 1  
919

920 It has been proved in Theorem 1 in (Haarnoja et al., 2017) that the constrained soft policy  $\pi_{\lambda;\theta} =$   
921  $\arg \max_{\pi} J_{r_{\theta}}(\pi) - \lambda J_c(\pi) + H(\pi)$  where we treat  $r_{\theta} - \lambda c$  as the new reward function. Recall  
922 that the dual function is  $G(\lambda; \theta) = \max_{\pi} J_{r_{\theta}}(\pi) - \lambda(J_c(\pi) - b) + H(\pi)$ , therefore, we know that  
923  $G(\lambda; \theta) = J_{r_{\theta}}(\pi_{\lambda;\theta}) - \lambda(J_c(\pi_{\lambda;\theta}) - b) + H(\pi_{\lambda;\theta})$ . Since  $\pi_{\lambda;\theta}$  is the optimal solution and the  $G$  is  
924 differentiable to the policy, we know that  $\frac{\partial G(\lambda;\theta)}{\partial \pi_{\lambda;\theta}(a|s)} = 0$  for every  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Note that  
925

$$926 \frac{\partial G(\lambda; \theta)}{\partial \pi_{\lambda;\theta}(a|s)} = \frac{\partial}{\partial \pi_{\lambda;\theta}(a|s)} \psi^{\pi_{\lambda;\theta}}(s, a) \left[ r_{\theta}(s, a) - \lambda c(s, a) - \log \pi_{\lambda;\theta}(a|s) \right] = 0 \quad (10)$$

$$927 \nabla_{\lambda} G(\lambda; \theta) = \nabla_{\lambda} \left[ J_{r_{\theta}}(\pi_{\lambda;\theta}) - \lambda J_c(\pi_{\lambda;\theta}) + H(\pi_{\lambda;\theta}) \right] - (J_c(\pi_{\lambda;\theta}) - b),$$

$$928 = \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \nabla_{\lambda} \left\{ \psi^{\pi_{\lambda;\theta}}(s, a) \left[ r_{\theta}(s, a) - \lambda c(s, a) - \log \pi_{\lambda;\theta}(a|s) \right] \right\} da ds - (J_c(\pi_{\lambda;\theta}) - b),$$

$$929 = \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \nabla_{\lambda} \pi_{\lambda;\theta}(a|s) \cdot \frac{\partial}{\partial \pi_{\lambda;\theta}(a|s)} \left\{ \psi^{\pi_{\lambda;\theta}}(s, a) \left[ r_{\theta}(s, a) - \lambda c(s, a) - \log \pi_{\lambda;\theta}(a|s) \right] \right\} da ds$$

$$930 - (J_c(\pi_{\lambda;\theta}) - b),$$

$$931 \stackrel{(a)}{=} b - J_c(\pi_{\lambda;\theta}),$$

932 where (a) follows (10).  
933

934 D.2 PROOF OF THEOREM 1 AND THEOREM 2

935 We first prove Theorem 2 and then prove Theorem 1. The fundamental logic is that we first prove  
936 that the dual problem has a unique optimal solution  $\lambda^*(\theta)$  and then we prove that the primal problem  
937 (i.e. the lower-level problem in (2)) and the dual problem have the same set of optimal solutions.  
938 Since the optimal solution of the dual problem is unique, then the optimal solution of the primal  
939 problem is also unique.

940 To show the optimal solution of the dual problem is unique, we prove that the dual function is strictly  
941 convex by showing that the Hessian of the dual function to  $\lambda$  is positive definite.

942 From Lemma 1, we know that  $\nabla_{\lambda} G(\lambda; \theta) = b - J_c(\pi_{\lambda;\theta})$ , therefore, we have that

$$943 \nabla_{\lambda\lambda}^2 G(\lambda; \theta) = -\nabla_{\lambda} J_c(\pi_{\lambda;\theta}),$$

$$944 = -\nabla_{\lambda} \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\lambda;\theta}(a_0|s_0) \left[ c(s_0, a_0) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) Q_c^{\pi_{\lambda;\theta}}(s_1) ds_1 \right] da_0 ds_0,$$

$$945 = -\int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \left\{ \nabla_{\lambda} \pi_{\lambda;\theta}(a_0|s_0) \cdot \left[ c(s_0, a_0) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) Q_c^{\pi_{\lambda;\theta}}(s_1) ds_1 \right] \right.$$

$$946 \left. + \pi_{\lambda;\theta}(a_0|s_0) \cdot \left[ \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) \nabla_{\lambda} Q_c^{\pi_{\lambda;\theta}}(s_1) ds_1 \right] \right\} da_0 ds_0,$$

$$947 = -\int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \left\{ \nabla_{\lambda} \pi_{\lambda;\theta}(a_0|s_0) \cdot Q_c^{\pi_{\lambda;\theta}}(s_0, a_0) \right.$$

$$948 \left. + \pi_{\lambda;\theta}(a_0|s_0) \cdot \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) \nabla_{\lambda} \left[ \int_{a_1 \in \mathcal{A}} \pi_{\lambda;\theta}(a_1|s_1) \cdot Q_c^{\pi_{\lambda;\theta}}(s_1, a_1) da_1 ds_1 \right] \right\} da_0 ds_0.$$

949 Keep the expansion, we can get that

$$950 \nabla_{\lambda\lambda}^2 G(\lambda; \theta) = -\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda;\theta}}(s) \nabla_{\lambda} \pi_{\lambda;\theta}(a|s) \cdot Q_c^{\pi_{\lambda;\theta}}(s, a) da ds,$$

$$951 = -\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda;\theta}}(s) \pi_{\lambda;\theta}(a|s) \nabla_{\lambda} \log \pi_{\lambda;\theta}(a|s) \cdot Q_c^{\pi_{\lambda;\theta}}(s, a) da ds,$$

$$\stackrel{(a)}{=} \int_{s \in \mathcal{S}} \psi^{\pi_{\lambda; \theta}}(s) \int_{a \in \mathcal{A}} \pi_{\lambda; \theta}(a|s) \left[ Q_c^{\pi_{\lambda; \theta}}(s, a) - V_c^{\pi_{\lambda; \theta}}(s) \right] Q_c^{\pi_{\lambda; \theta}}(s, a) da ds, \quad (11)$$

where (a) follows Lemma 5. Note that  $V_c^{\pi_{\lambda; \theta}}(s) = E_{a \sim \pi_{\lambda; \theta}(\cdot|s)}[Q_c^{\pi_{\lambda; \theta}}(s, a)]$ . If we use the random variable  $X_{sa}$  to denote  $Q_c^{\pi_{\lambda; \theta}}(s, a)$ , then its expectation is  $E_{a \sim \pi_{\lambda; \theta}(\cdot|s)}(X_{sa}) = V_c^{\pi_{\lambda; \theta}}(s)$ . We know that the variance  $Var(X_{sa}) = E[X_{sa}[X_{sa} - E(X_{sa})]]$ . Therefore, we can see that the equation (11) is actually a variance:

$$\nabla_{\lambda\lambda}^2 G(\lambda; \theta) = \int_{s \in \mathcal{S}} \psi^{\pi_{\lambda; \theta}}(s) Var(X_{sa}) ds.$$

From the expression (3), we know that  $\pi_{\lambda; \theta}$  is always stochastic. Therefore, the variance  $Var(X_{sa}) > 0$ . Then we know that  $\nabla_{\lambda\lambda}^2 G(\lambda; \theta) > 0$  and thus  $G$  is strictly convex. Therefore, the optimal solution  $\lambda^*(\theta)$  is unique.

Since  $G(\lambda; \theta)$  attains its minimum at  $\lambda^*(\theta)$ , the gradient of  $G$  at  $\lambda^*(\theta)$  should be zero, i.e.,  $J_c(\pi_{\lambda^*(\theta); \theta}) - b = 0$ . Let  $p^*$  and  $d^*$  be the optimal value of the primal problem and the dual problem. Since  $G(\lambda; \theta) = \max_{\pi} J_{r_{\theta}}(\pi) + H(\pi) - \lambda(J_c(\pi) - b)$ , we know that  $G(\lambda; \theta) \geq J_{r_{\theta}}(\pi) + H(\pi)$  for any  $(\lambda, \theta)$ , which means that  $d^* \geq p^*$ . Therefore, we have that:

$$p^* \leq d^* = G(\lambda^*(\theta); \theta) \stackrel{(b)}{=} J_{r_{\theta}}(\pi_{\lambda^*(\theta); \theta}) + H(\pi_{\lambda^*(\theta); \theta}) \leq p^*,$$

where (b) follows the fact that  $J_c(\pi_{\lambda^*(\theta); \theta}) - b = 0$ . Therefore,  $\pi_{\lambda^*(\theta); \theta}$  is an optimal solution of the primal problem. Suppose the primal problem has another optimal solution  $\pi'$ , then it holds that  $\pi' \in \arg \max_{\pi} G(\lambda^*(\theta); \theta)$ . However, it has been proved in Lemma 1 in (Zhou et al., 2017) that given an arbitrary  $\lambda$ , the optimal policy of  $\max_{\pi} J_{r_{\theta}}(\pi) + H(\pi) - \lambda(J_c(\pi) - b)$  is unique. Therefore,  $\pi_{\lambda^*(\theta); \theta}$  is the unique optimal solution of the primal problem (i.e., the lower-level problem in (2)).

### D.3 PROOF OF LEMMA 2

Since  $\lambda^*(\theta) = \arg \min G(\lambda; \theta)$ , we know that  $\nabla_{\lambda} G(\lambda^*(\theta); \theta) = 0$ . Therefore, we have that

$$\begin{aligned} \frac{d\nabla_{\lambda} G(\lambda^*(\theta); \theta)}{d\theta} &= 0, \\ \Rightarrow \nabla_{\theta\lambda}^2 G(\lambda^*(\theta); \theta) + \nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta) \nabla \lambda^*(\theta) &= 0, \\ \Rightarrow \nabla \lambda^*(\theta) &= -[\nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta)]^{-1} \nabla_{\theta\lambda}^2 G(\lambda^*(\theta); \theta). \end{aligned} \quad (12)$$

Now we take a look at the term  $\nabla_{\theta\lambda}^2 G(\lambda; \theta)$ . From Lemma 1, we know that  $\nabla_{\lambda} G(\lambda; \theta) = b - J_c(\pi_{\lambda; \theta})$ , therefore, we have that

$$\begin{aligned} \nabla_{\lambda\theta}^2 G(\lambda; \theta) &= -\nabla_{\theta} J_c(\pi_{\lambda; \theta}), \\ &= -\nabla_{\theta} \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\omega; \theta}(a_0|s_0) Q_c^{\pi_{\lambda; \theta}}(s_0, a_0) da_0 ds_0, \\ &= -\int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\omega; \theta}(a_0|s_0) \left[ \nabla_{\theta} \log \pi_{\omega; \theta}(a_0|s_0) \cdot Q_c^{\pi_{\lambda; \theta}}(s_0, a_0) + \nabla_{\theta} Q_c^{\pi_{\lambda; \theta}}(s_0, a_0) \right] da_0 ds_0, \\ &= -\int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\omega; \theta}(a_0|s_0) \left[ \nabla_{\theta} \log \pi_{\omega; \theta}(a_0|s_0) \cdot Q_c^{\pi_{\lambda; \theta}}(s_0, a_0) \right. \\ &\quad \left. - \nabla_{\theta} [c(s_0, a_0) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) Q_c^{\pi_{\lambda; \theta}}(s_1) ds_1] \right] da_0 ds_0, \\ &= -\int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\omega; \theta}(a_0|s_0) \left[ \nabla_{\theta} \log \pi_{\omega; \theta}(a_0|s_0) \cdot Q_c^{\pi_{\lambda; \theta}}(s_0, a_0) \right. \\ &\quad \left. - \gamma \int_{s_1 \in \mathcal{S}} P(s_1|s_0, a_0) \nabla_{\theta} \int_{a_1 \in \mathcal{A}} \pi_{\lambda; \theta}(a_1|s_1) Q_c^{\pi_{\lambda; \theta}}(s_1, a_1) da_1 ds_1 \right] da_0 ds_0. \end{aligned}$$

Keep the expansion, we can get that

$$\nabla_{\lambda\theta}^2 G(\lambda; \theta) = -\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s, a) \nabla_{\theta} \log \pi_{\lambda; \theta}(a|s) Q_c^{\pi_{\lambda; \theta}}(s, a) da ds,$$

$$\stackrel{(a)}{=} - \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s, a) \left[ Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s) \right] Q_c^{\pi_{\lambda; \theta}}(s, a) dads, \quad (13)$$

where (a) follows Lemma 5.

Now, we take the full gradient of  $\log \pi_{\lambda^*(\theta); \theta}(a|s)$  to  $\theta$ :

$$\begin{aligned} \frac{d \log \pi_{\lambda^*(\theta); \theta}(a|s)}{d\theta} &= \nabla_{\theta} \log \pi_{\lambda^*(\theta); \theta}(a|s) + \nabla_{\lambda} \log \pi_{\lambda^*(\theta); \theta}(a|s) \cdot \nabla \lambda^*(\theta), \\ &\stackrel{(b)}{=} Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) + (Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s)) [\nabla_{\lambda}^2 G(\lambda^*(\theta); \theta)]^{-1} \nabla_{\theta}^2 G(\lambda^*(\theta); \theta), \\ &\stackrel{(c)}{=} Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) - (Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s)). \\ &\frac{\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a) \left[ Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) \right] Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) dads}{\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a) \left[ Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s) \right] Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) dads}, \\ &= Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) - C_{\pi_{\lambda^*(\theta); \theta}}(Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s)), \end{aligned} \quad (14)$$

$$\text{where } C_{\pi_{\lambda; \theta}} \triangleq \frac{\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s, a) \left[ Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s) \right] Q_c^{\pi_{\lambda; \theta}}(s, a) dads}{\int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s, a) \left[ Q_c^{\pi_{\lambda; \theta}}(s, a) - V_c^{\pi_{\lambda; \theta}}(s) \right] Q_c^{\pi_{\lambda; \theta}}(s, a) dads}, \quad (b) \text{ follows Lemma 5, and}$$

(c) follows (11) and (13). Note that we can equivalently reformulate  $C_{\pi_{\lambda; \theta}}$  as:

$$C_{\pi_{\lambda; \theta}} = \frac{E_{(s, a) \sim \psi^{\pi_{\lambda; \theta}}(\cdot, \cdot)} \left[ (Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s)) Q_c^{\pi_{\lambda; \theta}}(s, a) \right]}{E_{(s, a) \sim \psi^{\pi_{\lambda; \theta}}(\cdot, \cdot)} \left[ (Q_c^{\pi_{\lambda; \theta}}(s, a) - V_c^{\pi_{\lambda; \theta}}(s)) Q_c^{\pi_{\lambda; \theta}}(s, a) \right]}. \quad (15)$$

Therefore, we can compute the hyper-gradient as:

$$\begin{aligned} \frac{dJ_r(\pi_{\lambda^*(\theta); \theta})}{d\theta} &\stackrel{(d)}{=} E_{(s, a) \sim \psi^{\pi_{\lambda^*(\theta); \theta}}} \left[ \frac{d \log \pi_{\lambda^*(\theta); \theta}}{d\theta} Q_r^{\pi_{\lambda^*(\theta); \theta}}(s, a) \right], \\ &\stackrel{(e)}{=} E_{(s, a) \sim \psi^{\pi_{\lambda^*(\theta); \theta}}} \left[ \left( Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) - C_{\pi_{\lambda^*(\theta); \theta}}(Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s)) \right) \right. \\ &\quad \left. Q_r^{\pi_{\lambda^*(\theta); \theta}}(s, a) \right], \end{aligned}$$

where (d) follows the standard result of policy gradient (Sutton & Barto, 2018), and (e) follows (14).

#### D.4 PROOF OF LEMMA 3

Since  $\pi_{\tilde{n}+1}(a|s) \propto \exp(Q_{\lambda; \theta}^{\pi_{\tilde{n}}}(s, a))$ , from Appendix C, we can see that  $\pi_{\tilde{n}+1}(a|s) = \frac{\exp(Q_{\lambda; \theta}^{\pi_{\tilde{n}}}(s, a))}{\exp(V_{\lambda; \theta}^{\pi_{\tilde{n}}}(s))}$ .

$$\begin{aligned} |\log \pi_{\tilde{n}+1}(a|s) - \log \pi_{\lambda; \theta}(a|s)| &= |Q_{\lambda; \theta}^{\pi_{\tilde{n}}}(s, a) - V_{\lambda; \theta}^{\pi_{\tilde{n}}}(s) - Q_{\lambda; \theta}(s, a) + V_{\lambda; \theta}(s)|, \\ &\stackrel{(a)}{=} |Q_{\lambda; \theta}(s, a) - Q_{\lambda; \theta}^{\pi_{\tilde{n}}}(s, a) + V_{\lambda; \theta}(s) - V_{\lambda; \theta}^{\pi_{\tilde{n}}}(s)|, \\ &\stackrel{(b)}{\leq} |Q_{\lambda; \theta}(s, a) - \mathcal{T}_{\lambda; \theta}^{\text{soft}}(Q_{\lambda; \theta}^{\pi_{\tilde{n}-1}})(s, a) + V_{\lambda; \theta}(s) - \mathcal{T}_{\lambda; \theta}^{\text{soft}}(V_{\lambda; \theta}^{\pi_{\tilde{n}-1}})(s)|, \\ &\stackrel{(c)}{\leq} |\mathcal{T}_{\lambda; \theta}^{\text{soft}}(Q_{\lambda; \theta})(s, a) - \mathcal{T}_{\lambda; \theta}^{\text{soft}}(Q_{\lambda; \theta}^{\pi_{\tilde{n}-1}})(s, a) + \mathcal{T}_{\lambda; \theta}^{\text{soft}}(V_{\lambda; \theta})(s) - \mathcal{T}_{\lambda; \theta}^{\text{soft}}(V_{\lambda; \theta}^{\pi_{\tilde{n}-1}})(s)|, \\ &\stackrel{(d)}{\leq} \gamma \left[ |Q_{\lambda; \theta}(s, a) - Q_{\lambda; \theta}^{\pi_{\tilde{n}-1}}(s, a) + V_{\lambda; \theta}(s) - V_{\lambda; \theta}^{\pi_{\tilde{n}-1}}(s)| \right], \\ &\stackrel{(e)}{\leq} \gamma^{\tilde{n}+1} \left[ |Q_{\lambda; \theta}(s, a) - Q_{\lambda; \theta}^{\pi_0}(s, a) + V_{\lambda; \theta}(s) - V_{\lambda; \theta}^{\pi_0}(s)| \right], \end{aligned}$$

where (a) follows policy improvement theorem (Theorem 4 in (Haarnoja et al., 2017)) (note that  $Q_{\lambda;\theta}$  and  $V_{\lambda;\theta}$  are the optimal Q/value functions under  $(\lambda, \theta)$ ), (b) follows Lemma 7, (c) follows the fact that the optimal Q/value functions are the fixed points of the contraction operator  $\mathcal{T}_{\lambda;\theta}^{\text{soft}}$ , and (d) follows Lemma 6.

**Lemma 8.** For any  $\theta$  and  $\bar{n} \geq 0$ , it holds that  $\nabla_{\lambda\lambda}^2 G(\lambda_{\bar{n}}; \theta) \succeq \tau_G I$  where  $\tau_G$  is a positive constant.

*Proof.* It has been proved in Subsection D.2 that  $\nabla_{\lambda\lambda}^2 G(\lambda; \theta) \succ 0$ . To prove that  $\nabla_{\lambda\lambda}^2 G(\lambda_{\bar{n}}; \theta) \succeq \tau_G I$ , we first prove that  $\nabla_{\lambda\lambda}^2 G(\lambda; \theta)$  is continuous in  $\lambda$  and then prove that the trajectory of  $\lambda_{\bar{n}}$  is bounded within a compact set for any  $\bar{n} \geq 0$ .

From (11), we know that

$$\nabla_{\lambda\lambda}^2 G(\lambda; \theta) = E^{\pi_{\lambda;\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t \left( Q_c^{\pi_{\lambda;\theta}}(s, a) - V_c^{\pi_{\lambda;\theta}}(s, a) \right) Q_c^{\pi_{\lambda;\theta}}(s, a) \right].$$

Since  $\pi_{\lambda;\theta}$ ,  $Q_c^{\pi_{\lambda;\theta}}(s, a)$ , and  $V_c^{\pi_{\lambda;\theta}}(s, a)$  are differentiable to  $\lambda$ , we know that  $\nabla_{\lambda\lambda}^2 G(\lambda; \theta)$  is continuous to  $\lambda$ . Now, we show that the trajectory of  $\lambda_{\bar{n}}$  is bounded within a compact set.

$$\begin{aligned} & \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2 = \|\lambda_{\bar{n}} - \alpha_{\bar{n}} g_{\lambda_{\bar{n}}; \theta} - \lambda^*(\theta)\|^2, \\ & = \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 - \alpha_{\bar{n}} \langle g_{\lambda_{\bar{n}}; \theta}, \lambda_{\bar{n}} - \lambda^*(\theta) \rangle, \\ & = \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 - \alpha_{\bar{n}} \langle \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta), \lambda_{\bar{n}} - \lambda^*(\theta) \rangle \\ & \quad - \alpha_{\bar{n}} \langle g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta), \lambda_{\bar{n}} - \lambda^*(\theta) \rangle, \\ & \leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 - \alpha_{\bar{n}} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)] \\ & \quad - \alpha_{\bar{n}} \langle g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta), \lambda_{\bar{n}} - \lambda^*(\theta) \rangle, \tag{16} \\ & \leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 - \alpha_{\bar{n}} \langle g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta), \lambda_{\bar{n}} - \lambda^*(\theta) \rangle, \\ & \leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} \|g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta)\| \cdot \|\lambda_{\bar{n}} - \lambda^*(\theta)\|, \\ & \stackrel{(a)}{\leq} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|, \\ & = \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| - \sum_{i=0}^{\bar{n}-1} \alpha_i g_{\lambda_i; \theta}, \\ & \leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| + \alpha_{\bar{n}} C \gamma^{\bar{n}} \left\| \sum_{i=0}^{\bar{n}-1} \alpha_i g_{\lambda_i; \theta} \right\|, \\ & \stackrel{(b)}{\leq} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \left( b + \frac{c_{\max}}{1-\gamma} \right)^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| + \alpha_{\bar{n}} C \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i \left( b + \frac{c_{\max}}{1-\gamma} \right), \tag{17} \end{aligned}$$

where (a) follows (18), and (b) follows that  $\|g_{\lambda; \theta}\| = \|b - J_c(\pi_{\lambda; \theta})\| \leq b + \frac{c_{\max}}{1-\gamma}$ . Now we show that  $\sum_{\bar{n}=1}^{\infty} \alpha_{\bar{n}} \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i$  is bounded. Since  $\alpha_i \propto 1/i^{\bar{n}}$ , we know that  $\sum_{i=0}^{\bar{n}-1} \alpha_i = O(\bar{n}^{1-\bar{n}})$ . Therefore, we know that  $\alpha_{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i = O(\bar{n}^{1-2\bar{n}}) \leq \bar{C}$  where  $\bar{C}$  is a positive constant. Therefore,  $\sum_{\bar{n}=1}^{\infty} \alpha_{\bar{n}} \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i \leq \bar{C} \sum_{\bar{n}=1}^{\infty} \gamma^{\bar{n}}$  is bounded. Now, we sum the both sides of (17) from  $\bar{n} = 1$  to  $\bar{N} - 1$ :

$$\begin{aligned} & \sum_{\bar{n}=0}^{\bar{N}-1} \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2, \\ & \leq \sum_{\bar{n}=0}^{\bar{N}-1} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \left( b + \frac{c_{\max}}{1-\gamma} \right)^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| + \alpha_{\bar{n}} C \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i \left( b + \frac{c_{\max}}{1-\gamma} \right), \\ & \Rightarrow \|\lambda_{\bar{N}} - \lambda^*(\theta)\|^2, \\ & \leq \|\lambda_0 - \lambda^*(\theta)\|^2 + \sum_{\bar{n}=0}^{\bar{N}-1} \alpha_{\bar{n}}^2 \left( b + \frac{c_{\max}}{1-\gamma} \right)^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| + \alpha_{\bar{n}} C \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i \left( b + \frac{c_{\max}}{1-\gamma} \right), \end{aligned}$$

$$\leq \|\lambda_0 - \lambda^*(\theta)\|^2 + \sum_{\bar{n}=0}^{\infty} \alpha_{\bar{n}}^2 (b + \frac{c_{\max}}{1-\gamma})^2 + \alpha_{\bar{n}} C \gamma^{\bar{n}} \|\lambda_0 - \lambda^*(\theta)\| + \alpha_{\bar{n}} C \gamma^{\bar{n}} \sum_{i=0}^{\bar{n}-1} \alpha_i (b + \frac{c_{\max}}{1-\gamma}).$$

Note that  $\alpha_{\bar{n}} \propto \frac{1}{(\bar{n}+1)^{\bar{\eta}}}$  and  $\bar{\eta} \in (\frac{1}{2}, 1)$ , it is obvious that  $|\lambda_{\bar{N}} - \lambda^*(\theta)|^2$  is bounded. Therefore, the trajectory of  $\lambda_{\bar{n}}$  is bounded for any  $\bar{n} \geq 0$ . Therefore, we can always find a positive constant  $\tau_G$  such that  $\nabla_{\lambda}^2 G(\lambda; \theta) \succeq \tau_G I$ .  $\square$

#### D.5 PROOF OF LEMMA 4

We first quantify the gradient approximation error  $|\nabla_{\lambda} G(\lambda; \theta) - g_{\lambda; \theta}|$  and then show the convergence of the middle loop.

$$\begin{aligned} |\nabla_{\lambda} G(\lambda; \theta) - g_{\lambda; \theta}| &= |J_c(\pi_{\lambda; \theta}) - J_c(\hat{\pi}_{\lambda; \theta})|, \\ &= \left| \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} [\psi^{\pi_{\lambda; \theta}}(s, a) - \psi^{\hat{\pi}_{\lambda; \theta}}(s, a)] c(s, a) da ds \right|, \\ &\leq c_{\max} \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} |\psi^{\pi_{\lambda; \theta}}(s, a) - \psi^{\hat{\pi}_{\lambda; \theta}}(s, a)| da ds, \\ &\stackrel{(a)}{\leq} c_{\max} C_d \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} |Q_{\lambda; \theta}(s, a) - Q_{\hat{\lambda}; \theta}^{\hat{\pi}_{\lambda; \theta}}(s, a)| da ds, \\ &\leq c_{\max} C_d C_{SA} \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \{|Q_{\lambda; \theta}(s, a) - Q_{\hat{\lambda}; \theta}^{\hat{\pi}_{\lambda; \theta}}(s, a)|\}, \\ &\stackrel{(b)}{\leq} c_{\max} C_d C_{SA} \gamma^{\bar{N}_{\bar{n}}} \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \{|Q_{\lambda; \theta}(s, a) - Q_{\hat{\lambda}; \theta}^{\pi_0}(s, a)|\}, \\ &= C \gamma^{\bar{N}_{\bar{n}}} \end{aligned} \tag{18}$$

where (a) follows step (iv) of equation (64) in (Zeng et al., 2022) and  $C_d$  is a positive constant,  $C_{SA}$  can be any positive constant that is larger than the area of  $\mathcal{S} \times \mathcal{A}$ , (b) follows the proof in Subsection D.4, and  $C = c_{\max} C_d C_{SA} \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \{|Q_{\lambda; \theta}(s, a) - Q_{\hat{\lambda}; \theta}^{\pi_0}(s, a)|\}$ .

Now, we quantify the convergence of the middle loop. From the expression (11) of  $\nabla_{\lambda}^2 G(\lambda; \theta)$ , we know that  $\|\nabla_{\lambda}^2 G(\lambda; \theta)\| \leq \frac{2c_{\max}^2}{(1-\gamma)^2}$ . From (16), we know that:

$$\begin{aligned} \alpha_{\bar{n}} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)] &\leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2 \\ &+ \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 - \alpha_{\bar{n}} \langle g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta), \lambda_{\bar{n}} - \lambda^*(\theta) \rangle, \\ &\leq \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} \|g_{\lambda_{\bar{n}}; \theta} - \nabla_{\lambda} G(\lambda_{\bar{n}}; \theta)\| \cdot \|\lambda_{\bar{n}} - \lambda^*(\theta)\|, \\ &\stackrel{(c)}{\leq} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 \|g_{\lambda_{\bar{n}}; \theta}\|^2 + \alpha_{\bar{n}} \gamma^{\bar{n}} \tilde{C}, \\ &\stackrel{(d)}{\leq} \|\lambda_{\bar{n}} - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{n}+1} - \lambda^*(\theta)\|^2 + \alpha_{\bar{n}}^2 (b + \frac{c_{\max}}{1-\gamma})^2 + \alpha_{\bar{n}} \gamma^{\bar{n}} \tilde{C}, \end{aligned} \tag{19}$$

where (c) follows (18) and the fact that  $\|\lambda_{\bar{n}} - \lambda^*(\theta)\|$  is bounded (proved in Lemma 8),  $\tilde{C}$  is a positive constant, and (d) follows that  $\|g_{\lambda; \theta}\| = \|b - J_c(\pi_{\lambda; \theta})\| \leq b + \frac{c_{\max}}{1-\gamma}$ . Telescoping (19) from  $\bar{n} = 0$  to  $\bar{N} - 1$ :

$$\begin{aligned} &\sum_{\bar{n}=0}^{\bar{N}-1} \alpha_{\bar{n}} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)], \\ &\leq \|\lambda_0 - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{N}} - \lambda^*(\theta)\|^2 + \sum_{\bar{n}=0}^{\bar{N}-1} \alpha_{\bar{n}}^2 (b + \frac{c_{\max}}{1-\gamma})^2 + \sum_{\bar{n}=0}^{\bar{N}-1} \alpha_{\bar{n}} \gamma^{\bar{n}} \tilde{C}. \end{aligned}$$

Since  $\alpha_{\bar{n}} = \frac{1}{(\bar{n}+1)^{\bar{\eta}}}$  and  $\bar{\eta} \in (\frac{1}{2}, 1)$ , there is a positive constant  $D_{\max}$  such that  $\sum_{\bar{n}=0}^{\bar{N}} \alpha_{\bar{n}}^2 (b + \frac{c_{\max}}{1-\gamma})^2 + \sum_{\bar{n}=0}^{\bar{N}} \alpha_{\bar{n}} \gamma^{\bar{n}-1} \tilde{C} \leq D_{\max}$ . Therefore, we have that

$$\sum_{\bar{n}=0}^{\bar{N}-1} \frac{1}{\bar{N}^{\bar{\eta}}} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)] \leq \sum_{\bar{n}=0}^{\bar{N}-1} \alpha_{\bar{n}} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)],$$

$$\begin{aligned}
&\leq \|\lambda_0 - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{N}} - \lambda^*(\theta)\|^2 + D_{\max}, \\
&\Rightarrow \frac{1}{\bar{N}} \sum_{\bar{n}=0}^{\bar{N}-1} [G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta)] \leq \frac{1}{\bar{N}^{1-\bar{\eta}}} \left[ \|\lambda_0 - \lambda^*(\theta)\|^2 - \|\lambda_{\bar{N}} - \lambda^*(\theta)\|^2 + D_{\max} \right].
\end{aligned} \tag{20}$$

Therefore, we have that

$$\begin{aligned}
\|\hat{\lambda}(\theta) - \lambda^*(\theta)\| &\stackrel{(e)}{\leq} \frac{2}{\tau_G} [G(\hat{\lambda}(\theta); \theta) - G(\lambda^*(\theta); \theta)] \stackrel{(f)}{\leq} \frac{2}{\tau_G} \left[ \frac{1}{\bar{N}} \sum_{\bar{n}=0}^{\bar{N}-1} G(\lambda_{\bar{n}}; \theta) - G(\lambda^*(\theta); \theta) \right], \\
&\stackrel{(g)}{\leq} O\left(\frac{1}{\bar{N}^{1-\bar{\eta}}}\right),
\end{aligned} \tag{21}$$

where (e) follows the fact that  $G(\lambda; \theta)$  is  $\tau_G$ -strongly convex (Lemma 8), (f) follows Jensen's inequality (note that  $\hat{\lambda}(\theta) = \frac{1}{\bar{N}} \sum_{\bar{n}=0}^{\bar{N}-1} \lambda_{\bar{n}}$ ), and (g) follows (20).

Now, we take a look at the term

$$\begin{aligned}
&|\log \pi_{\lambda^*(\theta); \theta}(a|s) - \log \hat{\pi}_{\hat{\lambda}(\theta); \theta}(a|s)|, \\
&\leq |\log \pi_{\lambda^*(\theta); \theta}(a|s) - \log \pi_{\hat{\lambda}(\theta); \theta}(a|s)| + |\log \pi_{\hat{\lambda}(\theta); \theta}(a|s) - \log \hat{\pi}_{\hat{\lambda}(\theta); \theta}(a|s)|, \\
&\stackrel{(h)}{\leq} \frac{2c_{\max}}{1-\gamma} \|\hat{\lambda}(\theta) - \lambda^*(\theta)\| + |\log \pi_{\hat{\lambda}(\theta); \theta}(a|s) - \log \hat{\pi}_{\hat{\lambda}(\theta); \theta}(a|s)|, \\
&\stackrel{(i)}{\leq} O\left(\frac{1}{\bar{N}^{1-\bar{\eta}}} + \gamma^{\bar{N}}\right),
\end{aligned} \tag{22}$$

where (h) follows Lemma 5 such that  $|\nabla_{\lambda} \log \pi_{\lambda; \theta}| \leq \frac{2c_{\max}}{1-\gamma}$ , and (i) follows (21) and Lemma 3.

**Lemma 9.** *The upper-level loss function  $J_r(\pi_{\lambda^*(\theta); \theta})$  is  $L$ -Lipschitz and  $\bar{L}$ -smooth where  $L$  and  $\bar{L}$  are positive constants. Moreover, it holds that  $\|g_{\theta}\| \leq L$  and  $\|\nabla_{\theta} g_{\theta}\| \leq \bar{L}$ .*

*Proof.* This suffices to show that the norms  $\|\nabla J_r(\pi_{\lambda^*(\theta); \theta})\|$  and  $\|\nabla^2 J_r(\pi_{\lambda^*(\theta); \theta})\|$  are upper bounded by  $L$  and  $\bar{L}$ . From Subsection D.3, we know that

$$\begin{aligned}
&\frac{dJ_r(\pi_{\lambda^*(\theta); \theta})}{d\theta}, \\
&= E_{(s,a) \sim \psi^{\pi_{\lambda^*(\theta); \theta}}} \left[ \left( Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s) - C_{\pi_{\lambda^*(\theta); \theta}}(Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a) - V_c^{\pi_{\lambda^*(\theta); \theta}}(s)) \right) \right. \\
&\left. Q_r^{\pi_{\lambda^*(\theta); \theta}}(s, a) \right],
\end{aligned}$$

where  $C_{\pi_{\lambda^*(\theta); \theta}} = [\nabla_{\lambda \lambda}^2 G(\lambda^*(\theta); \theta)]^{-1} \nabla_{\lambda \theta}^2 G(\lambda^*(\theta); \theta)$ . Since  $\|[\nabla_{\lambda \lambda}^2 G(\lambda^*(\theta); \theta)]^{-1}\| \leq \frac{1}{\tau_G}$  (Lemma 8) and  $\|\nabla_{\lambda \theta}^2 G(\lambda; \theta)\| = \|E_{(s,a) \sim \psi^{\pi_{\lambda; \theta}}}[(Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s, a) - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s)) Q_c^{\pi_{\lambda; \theta}}(s, a)]\| \leq \frac{2C_2 c_{\max}}{(1-\gamma)^2}$ , we know that  $\|C_{\pi_{\lambda^*(\theta); \theta}}\| \leq \frac{2C_2 c_{\max}}{(1-\gamma)^2 \tau_G}$ . Therefore, it holds that

$$\left\| \frac{dJ_r(\pi_{\lambda^*(\theta); \theta})}{d\theta} \right\| \leq \frac{1}{1-\gamma} \cdot \left[ \left( \frac{2C_2}{1-\gamma} + \frac{2C_2 c_{\max}}{(1-\gamma)^2 \tau_G} \cdot \frac{2c_{\max}}{1-\gamma} \right) \frac{C_1}{1-\gamma} \right] = L. \tag{23}$$

Similarly, we can see that  $\|g_{\theta}\| \leq L$ .

Now, we take a look at the Hessian term  $\nabla^2 J_r(\pi_{\lambda^*(\theta); \theta})$ . We define  $h^{\pi}(s, a) \triangleq Q_{\nabla_{\theta} r_{\theta}}^{\pi}(s_t, a_t) - V_{\nabla_{\theta} r_{\theta}}^{\pi}(s_t) - C_{\pi}(Q_c^{\pi}(s_t, a_t) - V_c^{\pi}(s_t)) Q_r^{\pi}(s_t, a_t)$ ,  $H^{\pi}(s, a) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t h^{\pi}(s_t, a_t) | s_0 = s, a_0 = a]$ , and  $H^{\pi}(s) \triangleq E^{\pi}[\sum_{t=0}^{\infty} \gamma^t h^{\pi}(s_t, a_t) | s_0 = s]$ . We know that  $\|H^{\pi}(s, a)\| \leq L$  and  $\|H^{\pi}(s)\| \leq L$ . Therefore, we have that

$$\nabla^2 J_r(\pi_{\lambda^*(\theta); \theta}) = \nabla \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) H^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) da_0 ds_0,$$

$$\begin{aligned}
&= \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \left[ \nabla \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot H^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right. \\
&+ \left. \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot \nabla H^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right] da_0 ds_0, \\
&= \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{a_0 \in \mathcal{A}} \left[ \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \nabla \log \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot H^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right. \\
&+ \left. \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot \left( \nabla h^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1 | s_0, a_0) \nabla H^{\pi_{\lambda^*(\theta); \theta}}(s_1) ds_1 \right) \right] da_0 ds_0.
\end{aligned}$$

Keep the expansion, we know that

$$\begin{aligned}
&\nabla^2 J_r(\pi_{\lambda^*(\theta); \theta}), \\
&= \int_{(s, a) \in \mathcal{S} \times \mathcal{A}} \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a) \left[ \nabla h^{\pi_{\lambda^*(\theta); \theta}}(s, a) + \nabla \log \pi_{\lambda^*(\theta); \theta}(a | s) \cdot H^{\pi_{\lambda^*(\theta); \theta}}(s, a) \right] dads,
\end{aligned} \tag{24}$$

Since  $\nabla \log \pi_{\lambda^*(\theta); \theta}(a | s)$  and  $H^{\pi_{\lambda^*(\theta); \theta}}(s, a)$  are both bounded, the only thing left is to bound  $\|\nabla h^{\pi_{\lambda^*(\theta); \theta}}(s, a)\|$ . We aim to bound  $\|\nabla h^{\pi_{\lambda^*(\theta); \theta}}(s, a)\|$  by bounding each term in  $\|h^{\pi_{\lambda^*(\theta); \theta}}(s, a)\|$ . Note that  $\nabla C_{\pi_{\lambda^*(\theta); \theta}} = \nabla([\nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta)]^{-1} \nabla_{\lambda\theta}^2 G(\lambda^*(\theta); \theta)) = [\nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta)]^{-2} \nabla_{\lambda\theta}^2 G(\lambda^*(\theta); \theta) \frac{d}{d\theta}([\nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta)] + [\nabla_{\lambda\lambda}^2 G(\lambda^*(\theta); \theta)]^{-1} \frac{d}{d\theta} \nabla_{\lambda\theta}^2 G(\lambda^*(\theta); \theta))$ . Therefore, it suffices to show that  $\|\frac{d}{d\theta} Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a)\|$  and  $\|\frac{d}{d\theta} Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a)\|$  are bounded.

$$\begin{aligned}
\frac{d}{d\theta} Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a) &= \frac{d}{d\theta} E^{\pi_{\lambda^*(\theta); \theta}} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r_{\theta}(s_t | a_t) | s_0 = s, a_0 = a \right], \\
&= \frac{d}{d\theta} \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{s_0 \in \mathcal{A}} \left[ \nabla \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right. \\
&+ \left. \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot \nabla Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right] da_0 ds_0, \\
&= \frac{d}{d\theta} \int_{s_0 \in \mathcal{S}} P_0(s_0) \int_{s_0 \in \mathcal{A}} \left[ \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \nabla \log \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right. \\
&+ \left. \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot \left( \nabla_{\theta\theta}^2 r_{\theta}(s_0, a_0) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1 | s_0, a_0) \nabla V_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s_1) ds_1 \right) \right] da_0 ds_0.
\end{aligned}$$

Keep the expansion, we can get that

$$\begin{aligned}
&\frac{d}{d\theta} Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a), \\
&= \int_{(s, a) \in \mathcal{S} \times \mathcal{A}} \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a) \left[ \nabla_{\theta\theta}^2 r_{\theta}(s_0, a_0) + \nabla \log \pi_{\lambda^*(\theta); \theta}(a_0 | s_0) \cdot Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s_0, a_0) \right] dads.
\end{aligned}$$

Therefore, we can see that  $\|\frac{d}{d\theta} Q_{\nabla_{\theta r\theta}}^{\pi_{\lambda^*(\theta); \theta}}(s, a)\| \leq \frac{C_3}{1-\gamma} + \left[ \frac{2c_{\max}}{(1-\gamma)} \cdot \frac{C_2}{\tau_G(1-\gamma)} + \frac{2C_2}{(1-\gamma)} \right] \cdot \frac{C_2}{(1-\gamma)^2}$ .

Similarly, we can also see that  $\frac{d}{d\theta} Q_c^{\pi_{\lambda^*(\theta); \theta}}(s, a)$  is also bounded. Therefore, we can find a positive constant  $\bar{L}$  such that  $\|\nabla^2 J_r(\pi_{\lambda^*(\theta); \theta})\| \leq \bar{L}$ . With the same procedure, we can see that  $\|\nabla_{\theta} g_{\theta}\| \leq \bar{L}$ .  $\square$

**Lemma 10.** *The hyper-gradient approximation error is upper bounded, i.e.,  $\|\frac{d}{d\theta} J_r(\pi_{\lambda^*(\theta); \theta}) - g_{\theta}\| \leq O(\gamma^{\bar{N}} + \frac{1}{\bar{N}^{1-\bar{q}}})$ .*

*Proof.*

$$\begin{aligned}
&\left\| \frac{d}{d\theta} J_r(\pi_{\lambda^*(\theta); \theta}) - g_{\theta} \right\|, \\
&\leq \left\| \int_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left[ \psi^{\hat{\pi}_{\lambda}(\theta); \theta}(s, a) h^{\hat{\pi}_{\lambda}(\theta); \theta}(s, a) - \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a) h^{\pi_{\lambda^*(\theta); \theta}}(s, a) \right] dads \right\|,
\end{aligned}$$



$$\begin{aligned}
1296 & \stackrel{(a)}{\leq} (1-\gamma)L \int_{(s,a) \in \mathcal{S} \times \mathcal{A}} \|\psi^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a) - \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a)\| d\mu, \\
1297 & \\
1298 & \\
1299 & \leq (1-\gamma)L \int_{(s,a) \in \mathcal{S} \times \mathcal{A}} \|\psi^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a) - \psi^{\pi_{\hat{\lambda}(\theta); \theta}}(s, a)\| + \|\psi^{\pi_{\hat{\lambda}(\theta); \theta}}(s, a) - \psi^{\pi_{\lambda^*(\theta); \theta}}(s, a)\| d\mu, \\
1300 & \\
1301 & \stackrel{(b)}{\leq} (1-\gamma)C_d C_{SA} \left[ \max\{|Q_{\hat{\lambda}(\theta); \theta}(s, a) - Q_{\hat{\lambda}(\theta); \theta}^{\hat{\pi}_{\hat{\lambda}(\theta); \theta}}(s, a)|\} + \max\{|Q_{\hat{\lambda}(\theta); \theta}(s, a) - Q_{\lambda^*(\theta); \theta}(s, a)|\} \right], \\
1302 & \\
1303 & \stackrel{(c)}{\leq} (1-\gamma)C_d C_{SA} \left[ \gamma^{\bar{N}} \max\{|Q_{\hat{\lambda}(\theta); \theta}(s, a) - Q_{\hat{\lambda}(\theta); \theta}^{\pi_0}(s, a)|\} \right. \\
1304 & \\
1305 & \left. + \max\{|Q_{\hat{\lambda}(\theta); \theta}(s, a) - Q_{\lambda^*(\theta); \theta}(s, a)|\} \right], \\
1306 & \\
1307 & \stackrel{(d)}{\leq} (1-\gamma)C_d C_{SA} \left[ \gamma^{\bar{N}} \max\{|Q_{\hat{\lambda}(\theta); \theta}(s, a) - Q_{\hat{\lambda}(\theta); \theta}^{\pi_0}(s, a)|\} + \frac{c_{\max}}{1-\gamma} \|\lambda^*(\theta) - \hat{\lambda}(\theta)\| \right], \\
1308 & \\
1309 & \stackrel{(e)}{\leq} O(\gamma^{\bar{N}} + \frac{1}{\bar{N}^{1-\eta}}), \\
1310 & \\
1311 & \\
1312 &
\end{aligned}$$

1313 where (a) follows (23), (b) follows step (iv) of equation (64) in (Zeng et al., 2022), (c) follows  
1314 (18), (d) follows the fact that  $\|\nabla_{\lambda} Q_{\lambda; \theta}(s, a)\| = |Q_c^{\pi_{\lambda; \theta}}(s, a)| \leq \frac{c_{\max}}{1-\gamma}$ , and (e) follows Lemma  
1315 4.  $\square$

## 1316 D.6 PROOF OF THEOREM 3

1317 We define a function  $f(\theta)$  such that  $\nabla f(\theta) = g_{\theta}$  and  $\nabla^2 f(\theta) = \nabla_{\theta} g_{\theta}$ , therefore, we have that

$$\begin{aligned}
1318 & \\
1319 & \\
1320 & f(\theta_{n+1}) \leq f(\theta_n) + \langle \nabla f(\theta_n), \theta_{n+1} - \theta_n \rangle + \frac{\bar{L}}{2} \|\theta_{n+1} - \theta_n\|^2, \\
1321 & \\
1322 & = f(\theta_n) - \beta_n \|\nabla f(\theta_n)\|^2 + \frac{\bar{L}\beta_n^2}{2} \|\nabla f(\theta_n)\|^2, \\
1323 & \\
1324 & \Rightarrow \beta_n \|\nabla f(\theta_n)\|^2 \leq f(\theta_n) - f(\theta_{n+1}) + \frac{\bar{L}\beta_n^2}{2} \|\nabla f(\theta_n)\|^2. \quad (25) \\
1325 &
\end{aligned}$$

1326 Telescoping (25) from  $n = 0$  to  $N - 1$ , we have that

$$\begin{aligned}
1327 & \\
1328 & \sum_{n=0}^{N-1} \beta_n \|\nabla f(\theta_n)\|^2 \leq f(\theta_0) - f(\theta_N) + \sum_{n=0}^{N-1} \frac{\bar{L}\beta_n^2}{2} \|\nabla f(\theta_n)\|^2, \\
1329 & \\
1330 & \Rightarrow \sum_{n=0}^{N-1} \frac{1}{N\eta} \|\nabla f(\theta_n)\|^2 \leq \sum_{n=0}^{N-1} \beta_n \|\nabla f(\theta_n)\|^2 \stackrel{(a)}{\leq} f(\theta_0) - f(\theta_N) + \sum_{n=0}^{\infty} \frac{L^2 \bar{L} \beta_n^2}{2}, \\
1331 & \\
1332 & \Rightarrow \frac{1}{N} \sum_{n=0}^{N-1} \|\nabla f(\theta_n)\|^2 \leq \frac{1}{N^{1-\eta}} \left( f(\theta_0) - f(\theta_N) + \sum_{n=0}^{\infty} \frac{L^2 \bar{L} \beta_n^2}{2} \right) \stackrel{(b)}{=} O\left(\frac{1}{N^{1-\eta}}\right), \quad (26) \\
1333 & \\
1334 & \\
1335 &
\end{aligned}$$

1336 where (a) follows Lemma 9, and (b) follows the fact that  $\sum_{n=0}^{\infty} \beta_n^2$  is bounded as  $\beta_n = \frac{1}{(n+1)^\eta}$  and  
1337  $\eta \in (\frac{1}{2}, 1)$ . Therefore, we have that

$$\begin{aligned}
1338 & \\
1339 & \frac{1}{N} \sum_{n=0}^{N-1} \|\nabla J_r(\pi_{\lambda^*(\theta); \theta})\|^2 \leq \frac{1}{N} \sum_{n=0}^{N-1} \left( \|\nabla f(\theta_n)\|^2 + \|\nabla f(\theta_n) - \nabla J_r(\pi_{\lambda^*(\theta); \theta})\|^2 \right), \\
1340 & \\
1341 & \\
1342 & \stackrel{(c)}{\leq} O\left(\frac{1}{N^{1-\eta}} + \gamma^{2\bar{N}} + \frac{1}{N^{2-2\eta}}\right), \\
1343 &
\end{aligned}$$

1344 where (c) follows (26) and Lemma 10.

## 1345 D.7 PROOF OF THEOREM 4

1346 We first define  $F(\theta) \triangleq J_r(\pi_{\lambda^*(\theta); \theta})$ . Theorem 10 in (Agarwal et al., 2021) shows that the policy  
1347 gradient method can achieve global optimality asymptotically under softmax policy parameteriza-  
1348 tion. The constrained soft policy  $\pi_{\lambda^*(\theta); \theta} = \lim_{\bar{N} \rightarrow \infty} \hat{\pi}_{\hat{\lambda}(\theta); \theta}$  can be regarded as a softmax policy  
1349

parameterized by  $Q_{\lambda^*(\theta);\theta}$  but the decision variable in our case is  $\theta$  instead of  $Q_{\lambda^*(\theta);\theta}$ . However, we can still build on Theorem 10 in (Agarwal et al., 2021) by building connections between  $\theta$  and  $Q_{\lambda^*(\theta);\theta}$ . In specific, in order to use the result of Theorem 10 in (Agarwal et al., 2021), we need to prove (i)  $F(\theta_n)$  is monotonically increasing, i.e.,  $F(\theta_{n+1}) \geq F(\theta_n)$  for any  $n \geq 0$ ; (ii)  $\frac{dF(\bar{\theta})}{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}} = 0$  if  $\frac{dF(\bar{\theta})}{d\theta} = 0$  where  $Q_{\lambda^*(\bar{\theta});\bar{\theta}}$  is a vector with the length  $|\mathcal{S}| \times |\mathcal{A}|$  whose components are  $\{Q_{\lambda^*(\bar{\theta});\bar{\theta}}(s, a)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ . Once proving these two, given that  $\lim_{N \rightarrow \infty} \lim_{\bar{N} \rightarrow \infty} \frac{dF(\theta_N)}{d\theta} = 0$  from Theorem 3, we can use Theorem 10 in (Agarwal et al., 2021) to prove Theorem 4.

Now, we first show that  $F(\theta_n)$  is monotonically increasing. This is a straightforward result of Theorem 10.15 in (Beck, 2017) if we choose  $\beta_n \leq \frac{1}{L}$ . Note that Theorem 10 in (Agarwal et al., 2021) requires  $\beta_n \leq \frac{(1-\gamma)^3}{8}$ , so that we can choose  $\beta_n \leq \min\{\frac{1}{L}, \frac{(1-\gamma)^3}{8}\}$ . We next show that  $\frac{dF(\bar{\theta})}{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}} = 0$  if  $\frac{dF(\bar{\theta})}{d\theta} = 0$ .

We know that  $\frac{dF(\bar{\theta})}{d\theta} = \frac{dF(\bar{\theta})}{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}} \cdot \frac{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}}{d\theta}$ , so that it suffices to show that  $\frac{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}(s, a)}{d\theta} \neq 0$  for any  $(s, a)$  and any  $\theta$ . Therefore, we have that

$$\begin{aligned} \frac{dQ_{\lambda^*(\bar{\theta});\bar{\theta}}(s, a)}{d\theta} &= \nabla_{\lambda} Q_{\lambda^*(\bar{\theta});\bar{\theta}}(s, a) \nabla \lambda^*(\bar{\theta}) + \nabla_{\theta} Q_{\lambda^*(\bar{\theta});\bar{\theta}}(s, a), \\ &\stackrel{(a)}{=} Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta});\bar{\theta}}}(s, a) - Q_c^{\pi_{\lambda^*(\bar{\theta});\bar{\theta}}}(s, a) \nabla \lambda^*(\bar{\theta}), \end{aligned} \quad (27)$$

where (a) follows (6) and (8). Now, we first prove that the term  $Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta});\bar{\theta}}}(s, a) = E^{\pi_{\lambda^*(\bar{\theta});\bar{\theta}}}[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r_{\bar{\theta}}(s_t, a_t)]$  is nonzero. Define  $l(\theta; \lambda, s, a) \triangleq E^{\pi_{\lambda; \theta}}[\sum_{t=0}^{\infty} \gamma^t r_{\theta}(s_t, a_t) | s_0 = s, a_0 = a]$ , and therefore  $\nabla_{\theta} l(\theta; \lambda, s, a) \triangleq E^{\pi_{\lambda; \theta}}[\sum_{t=0}^{\infty} \gamma^t r_{\theta}(s_t, a_t) | s_0 = s, a_0 = a]$ . We use  $\psi^{\pi}(s' | s, a)$  and  $\psi^{\pi}(s', a' | s, a)$  to denote the state and state-action visitation frequency when the initial state-action is  $(s, a)$ . Now, we take a look at the Hessian term  $\nabla_{\theta\theta}^2 l(\theta; \lambda, s, a)$ :

$$\begin{aligned} \nabla_{\theta\theta}^2 l(\theta; \lambda, s, a), \\ &= \nabla_{\theta\theta}^2 r_{\theta}(s, a) + \gamma \nabla_{\theta} \int_{s_1 \in \mathcal{S}} P(s_1 | s, a) \int_{a_1 \in \mathcal{A}} \pi_{\lambda; \theta}(a_1 | s_1) \nabla_{\theta} l(\theta; \lambda, s_1, a_1) da_1 ds_1, \\ &= \nabla_{\theta\theta}^2 r_{\theta}(s, a) + \gamma \int_{s_1 \in \mathcal{S}} P(s_1 | s, a) \int_{a_1 \in \mathcal{A}} \left[ \nabla_{\theta} \pi_{\lambda; \theta}(a_1 | s_1) \cdot \nabla_{\theta} l(\theta; \lambda, s_1, a_1) \right. \\ &\quad \left. + \pi_{\lambda; \theta}(a_1 | s_1) \cdot \nabla_{\theta\theta}^2 l(\theta; \lambda, s_1, a_1) \right] da_1 ds_1. \end{aligned}$$

Keep the expansion and note that  $\nabla_{\theta} l(\theta; \lambda, s, a) = Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s, a)$ , we can get that

$$\begin{aligned} \nabla_{\theta\theta}^2 l(\theta; \lambda, s, a) &= \int_{s' \in \mathcal{S}} \int_{a' \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s' | s, a) \nabla_{\theta} \pi_{\lambda; \theta}(a' | s') \cdot Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a') da' ds', \\ &= \int_{s' \in \mathcal{S}} \int_{a' \in \mathcal{A}} \psi^{\pi_{\lambda; \theta}}(s' | s, a) \pi_{\lambda; \theta}(a' | s') \nabla_{\theta} \log \pi_{\lambda; \theta}(a' | s') \cdot Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a') da' ds', \\ &\stackrel{(b)}{=} \int_{s' \in \mathcal{S}} \psi^{\pi_{\lambda; \theta}}(s' | s, a) \int_{a' \in \mathcal{A}} \pi_{\lambda; \theta}(a' | s') \left[ Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a') - V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s') \right] Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a') da' ds', \end{aligned} \quad (28)$$

where (b) follows Lemma 5. Note that  $V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s') = E_{a' \sim \pi_{\lambda; \theta}(\cdot | s')} [Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a')]$ . If we use the random variable  $Y_{s'a'}$  to denote  $Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s', a')$ , then its expectation is  $E_{a' \sim \pi_{\lambda; \theta}(\cdot | s')} (Y_{s'a'}) = V_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda; \theta}}(s')$ . We know that the variance  $Var(Y_{s'a'}) = E[Y_{s'a'} [Y_{s'a'} - E(Y_{s'a'})]]$ . Therefore, we can see that the equation (28) is actually a variance:

$$\nabla_{\theta\theta}^2 l(\theta; \lambda, s, a) = \int_{s' \in \mathcal{S}} \psi^{\pi_{\lambda; \theta}}(s' | s, a) Var(Y_{s'a'}) ds' \succeq 0.$$

Therefore, the function  $l(\theta; \lambda, s, a)$  is convex in  $\theta$  for any  $(\lambda, s, a)$ . If  $\nabla_{\theta} l(\bar{\theta}; \lambda, s, a) = 0$ , this means that  $l(\bar{\theta}; \lambda, s, a)$  achieves its optimum. However,  $l(\theta; \lambda, s, a)$  does not have an optimum, i.e.,  $l(\theta; \lambda, s, a)$  can be infinitely large or infinitely small because  $r_{\theta}(s, a)$  can be any arbitrarily large

value. This is a contradiction, therefore,  $l(\bar{\theta}; \lambda, s, a) \neq 0$  for any  $(\lambda, s, a)$ . Then,  $l(\bar{\theta}; \lambda^*(\bar{\theta}), s, a) = Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) \neq 0$ .

Recall from (27) that  $\frac{dQ_{\lambda^*(\bar{\theta}); \bar{\theta}}(s, a)}{d\theta} = Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) - Q_c^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) \nabla \lambda^*(\bar{\theta})$ . For any  $(s, a) \notin \mathcal{C}$ ,  $Q_c^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) = 0$  because the policy  $\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}$  satisfies the constraint of the lower-level problem in (1), i.e., avoiding the set  $\mathcal{C}$ . Therefore,  $\frac{dQ_{\lambda^*(\bar{\theta}); \bar{\theta}}(s, a)}{d\theta} = Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) \neq 0$ . For any  $(s, a) \in \mathcal{C}$ , we know that  $Q_c^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) = c(s, a)$  because the policy  $\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}$  avoids the set  $\mathcal{C}$  unless its starting state-action pair is in  $\mathcal{C}$ . Therefore, we can also design  $c(s, a)$  such that  $Q_{\nabla_{\theta} r_{\theta}}^{\pi_{\lambda^*(\bar{\theta}); \bar{\theta}}}(s, a) - c(s, a) \nabla \lambda^*(\bar{\theta}) \neq 0$  for  $(s, a) \in \mathcal{C}$ . Therefore, we can ensure that  $\frac{dQ_{\lambda^*(\bar{\theta}); \bar{\theta}}(s, a)}{d\theta} \neq 0$  for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Therefore,  $\frac{dF(\bar{\theta})}{dQ_{\lambda^*(\bar{\theta}); \bar{\theta}}} = 0$  if  $\frac{dF(\bar{\theta})}{d\theta} = 0$ .

## E RELATED WORKS

**XRL methods that has the potential to be used to improve RL performance.** There are some XRL methods that have the potential to improve the RL performance even if they do not mention that they can improve the RL performance. Value-max (Amir & Amir, 2018; Huang et al., 2018) use the value function  $V(s)$  to identify the states with highest value as critical points. We can perturb the actions on these critical states to improve the RL performance.

**Constrained reinforcement learning (CRL).** The lower-level problem in (1) is a CRL problem. The current works on CRL have two major categories: primal-dual approach and primal approach. The primal-dual approach (Achiam et al., 2017; Tessler et al., 2018; Stooke et al., 2020) converts the CRL problem into an unconstrained optimization problem by using the dual method. Our approach can be categorized as a primal-dual approach. The primal approach (Liu et al., 2020; Chow et al., 2018; Xu et al., 2021) enforce constraints via various designs of the objective function or the update process without an introduction of dual variables. However, these previous methods on CRL may not be suitable to the context of constrained bi-level optimization because they cannot guarantee that the upper-level problem in (1) is smooth after the lower-level problem is solved. The non-smoothness of the upper-level problem can make the constrained bi-level optimization problem difficult to solve. In contrast, our approach ensures the smoothness of the upper-level problem in (1) because we derive an analytical solution of the constrained soft policy and this policy is smooth w.r.t.  $\theta$ .

### E.1 COMPARISON WITH (HU ET AL., 2020)

Paper (Hu et al., 2020) studies how to utilize the domain knowledge to learn a shaping reward function and formulate a bi-level optimization problem. Here, we discussion our distinctions and improvements compared to (Hu et al., 2020) in terms of problem, assumption, algorithm, and theoretical analysis.

**Problem:** Our problem is not only to shape the reward function, but also discourage from visiting the set  $\mathcal{C}$ , therefore, we have a constrained bi-level optimization problem where the lower level is a constrained RL that is non-convex. Paper (Hu et al., 2020) only studies how to shape the reward so that its bi-level optimization problem is unconstrained, which is much easier to solve.

**Assumption:** Paper (Hu et al., 2020) assumes that the domain knowledge  $\hat{r}$  is given by humans while the domain knowledge  $\hat{r}$  is learned in our case. Moreover, paper (Hu et al., 2020) assumes that the shaping reward is linear, i.e., the shaping reward is  $r(s, a) + \theta(s, a)(r(s, a) - \hat{r}(s, a))$ . In contrast, our shaping reward class  $r_{\theta}(r(s, a), r(s, a) - \hat{r}(s, a))$  is more general and includes the linear shaping reward as a special case.

**Algorithm:** Paper (Hu et al., 2020) proposes several methods to compute the hyper-gradient (i.e., the gradient of the upper-level problem). However, they do not consider the practical issue, i.e., the lower-level problem cannot be fully solved in finite time. In contrast, our constrained bi-level optimization problem cannot be solved by the state-of-the-arts (Xu & Zhu, 2023; Khanduri et al., 2023) so that we develop a novel algorithm that solve the problem and we consider the practical issue, i.e., we cannot obtain the exact optimal solution in finite time.

**Theoretical analysis:** Paper (Hu et al., 2020) does not have theoretical guarantees at all. In contrast, we propose a systematic theoretical framework, which is one of our core contributions.

## F EXPERIMENT DETAILS

The code was running on a laptop whose CPU is Intel Core i9 12900k and GPU is NVIDIA RTX 3080. The operating system is Windows 10. We use a neural network to parameterize the learned reward function. The neural network has two hidden layers where each hidden layer has 64 neurons. The activation functions are respectively ReLU and Tanh.

**The delayed MuJoCo environments.** The delayed Mujoco environments are widely used in RL improvement literature (Zheng et al., 2018; Memarian et al., 2021; Oh et al., 2018) where the reward is accumulated by 20 time steps and only provided at the end. For example, for an episode with length 100 (i.e., 0, ..., 99), only the time steps 19, 39, 59, 79, 99 receive nonzero reward while all the other steps receive zero reward. The time step 19 receives the reward that is accumulated from time 0 to time 19. We can see that the delayed MuJoCo tasks have sparse reward.

### F.1 USING OTHER BASELINE RL ALGORITHMS

We provide the results of using PPO and TD3 as the baseline RL algorithm below:

Table 4: Using PPO as the baseline RL algorithm.

	Delayed HalfCheetah	Delayed Hopper	Delayed Walker	Delayed Ant	HalfCheetah	Hopper	Walker	Ant
PPO	179.97 ± 20.12	243.14 ± 25.87	175.36 ± 18.22	33.28 ± 5.10	221.58 ± 25.29	285.10 ± 36.29	203.19 ± 24.84	302.94 ± 28.29
UTILITY	422.83 ± 32.15	400.99 ± 21.45	395.15 ± 30.28	88.42 ± 5.54	457.03 ± 30.29	420.46 ± 35.94	402.40 ± 39.28	392.95 ± 39.10
RICE	226.54 ± 16.35	271.85 ± 18.73	205.15 ± 16.11	42.62 ± 3.48	238.17 ± 19.99	295.82 ± 21.18	229.52 ± 23.22	315.38 ± 22.38
SIL	232.55 ± 25.75	259.15 ± 22.43	183.27 ± 10.39	49.66 ± 3.21	271.92 ± 29.18	301.93 ± 28.37	284.44 ± 48.17	311.92 ± 24.73
LIR	341.38 ± 22.54	311.44 ± 21.18	315.77 ± 28.59	73.14 ± 6.28	319.29 ± 37.19	314.73 ± 30.18	326.28 ± 28.27	325.83 ± 18.37

Table 5: Using TD3 as the baseline RL algorithm.

	Delayed HalfCheetah	Delayed Hopper	Delayed Walker	Delayed Ant	HalfCheetah	Hopper	Walker	Ant
TD3	216.84 ± 18.50	252.84 ± 19.68	167.02 ± 21.52	43.15 ± 4.28	237.29 ± 22.51	283.75 ± 28.14	193.25 ± 22.15	364.82 ± 21.17
UTILITY	464.74 ± 27.06	422.13 ± 24.15	392.04 ± 24.77	105.24 ± 7.88	476.92 ± 27.15	433.15 ± 27.13	396.16 ± 28.59	483.11 ± 26.85
RICE	221.16 ± 16.88	269.17 ± 14.66	179.13 ± 12.83	46.09 ± 3.83	236.83 ± 19.26	289.99 ± 21.77	212.26 ± 21.04	371.65 ± 22.85
SIL	258.16 ± 23.92	308.63 ± 26.25	201.66 ± 12.95	58.62 ± 6.38	312.29 ± 24.88	315.17 ± 24.72	289.51 ± 31.43	336.16 ± 21.63
LIR	377.19 ± 26.56	346.10 ± 18.11	298.17 ± 23.77	88.19 ± 3.64	386.10 ± 33.83	366.28 ± 26.88	335.15 ± 20.62	381.17 ± 32.15

Table 4 and Table 5 show that UTILITY significantly outperforms the baselines when the baseline RL algorithms are PPO and TD3.

### F.2 EXPERIMENT RESULTS ON DENSE REWARD

To show that our method can also improve the performance on dense reward scenarios, we include the experiment results on the original MuJoCo environment (where the reward is dense) below:

Table 6: Experiment results (original MuJoCo environment with dense reward).

	SAC	UTILITY	RICE	SIL	LIR
HalfCheetah	686.40 ± 41.24	824.42 ± 42.18	701.44 ± 45.83	716.62 ± 39.17	718.25 ± 42.10
Hopper	238.14 ± 29.94	348.16 ± 26.32	242.99 ± 19.10	264.86 ± 27.11	277.58 ± 22.69
Walker2d	182.21 ± 24.14	269.14 ± 25.08	189.19 ± 39.11	197.63 ± 21.39	214.16 ± 21.24
Ant	299.79 ± 26.53	421.63 ± 25.16	322.26 ± 34.15	340.14 ± 26.53	351.14 ± 28.92

Table 6 shows the final results on the original MuJoCo environment (dense reward). We can observe that UTILITY achieves the highest reward and largely improves SAC.

### F.3 EXPERIMENT RESULTS WHEN THE EPISODE LENGTH IS 1000

We provide the experiment results when the episode length is 1,000.

Table 7: Results when the episode length is 1,000.

	Delayed HalfCheetah	Delayed Hopper	Delayed Walker	Delayed Ant	HalfCheetah	Hopper	Walker	Ant
SAC	3019.18 ± 88.69	1975.28 ± 66.10	1827.43 ± 108.19	3652.85 ± 59.70	5021.39 ± 102.15	3592.29 ± 88.62	3344.21 ± 95.58	5922.10 ± 128.19
UTILITY	4872.10 ± 86.39	3562.99 ± 74.44	3262.52 ± 79.01	4929.14 ± 74.26	5806.28 ± 69.01	4217.83 ± 75.62	3829.16 ± 73.07	6500.73 ± 96.18
RICE	3138.44 ± 35.66	2196.82 ± 64.29	1996.53 ± 50.82	3742.18 ± 49.45	5127.77 ± 42.72	3719.42 ± 63.17	3396.28 ± 50.27	5998.26 ± 62.16
SIL	3528.53 ± 56.03	2549.98 ± 85.43	2393.75 ± 56.56	3940.95 ± 60.10	5429.58 ± 74.95	3692.48 ± 35.84	3462.94 ± 44.82	6102.58 ± 57.29
LIR	3921.61 ± 142.06	2842.19 ± 71.02	2511.26 ± 38.51	4020.29 ± 61.02	5420.92 ± 59.29	3792.40 ± 101.93	3502.94 ± 24.64	6204.56 ± 39.12

Table 7 shows that UTILITY outperforms the baselines when the episode length is 1,000.

#### F.4 ABLATION STUDY

Since our method has two components to improve the performance: the shaping reward and the constrained formulated by the “misleading” state-action pairs. Here, we separately study the effect of the learned shaping reward and the constraint. In specific, we test the performance of the shaping only method and the constraint only method, and provide the results below:

Table 8: Ablation study for dense reward.

	SAC	UTILITY	shaping only	constraint only
HalfCheetah	686.40 ± 51.24	824.42 ± 42.18	764.25 ± 48.11	701.19 ± 47.43
Hopper	238.14 ± 29.94	348.16 ± 26.32	311.78 ± 34.24	268.15 ± 42.66
Walker2d	182.21 ± 24.14	269.14 ± 25.08	242.18 ± 29.62	196.77 ± 22.19
Ant	299.79 ± 26.53	421.63 ± 25.16	396.05 ± 29.18	317.12 ± 34.59

Table 8 shows that both the learned shaping reward and the constraint can improve the performance, and the shaping reward has a larger impact. Moreover, even if we only use the shaping reward, the performance is better than LIR. This is because our shaping reward uses the domain knowledge formulated by the high-level explanation. Even if we only use the constraint, the performance is better than RICE. The reason is that the “misleading” state-action pairs we find are the points that lead to the failure, and thus avoiding these state-action pairs can improve performance. In contrast, RICE finds the states that are most influential cumulative reward, however, these states may not be the states that lead the RL agent to be non-optimal.

#### F.5 FIDELITY OF THE GENERATED TWO-LEVEL EXPLANATION

The fidelity means the correctness of the two-level explanation (Guo et al., 2021b; Cheng et al., 2023). Since the two-level explanation is to explain why the RL agent (i.e., SAC) is not optimal, one way to validate the fidelity of the explanation is to see whether the performance improves after we improve from the explanations. From the last two columns in Table 8, we can see that both the high-level and low-level explanations are the correct explanations because both the shaping only method and the constraint only method improve the performance. Moreover, the shaping only method (the fourth column in Table 8) has a higher cumulative reward than LIR (the last column in Table 6), and the constraint only method (the last column in Table 8) has a higher cumulative reward than RICE (the fourth column in Table 6). This shows the high fidelity of our two-level explanation.

Table 9: Fidelity comparison for dense reward.

	SAC	shaping only (ours)	shaping without $r - \hat{r}$	constraint only (ours)	RICE+constraint
HalfCheetah	686.40 ± 51.24	764.25 ± 48.11	710.99 ± 35.16	701.19 ± 47.43	672.15 ± 25.16
Hopper	238.14 ± 29.94	311.78 ± 34.24	279.12 ± 27.44	268.15 ± 42.66	211.03 ± 21.20
Walker2d	182.21 ± 24.14	242.18 ± 29.62	213.15 ± 18.94	196.77 ± 22.19	177.19 ± 15.33
Ant	299.79 ± 26.53	396.05 ± 29.18	359.14 ± 22.85	317.12 ± 34.59	288.19 ± 18.02

Table 9 shows that both the high-level and low-level explanations of our method have higher fidelity. The method “RICE+constraint” has even worse performance than SAC because the critical states influential to the cumulative reward may be the states that lead to high cumulative reward, and thus constraining them may even make the performance worse. However, even if we do not constrain

these states but use the fine-tune method as in (Cheng et al., 2023) instead, our constraint-only method (the last column in Table 8) still outperforms RICE (the fourth column in Table 6). For the high-level explanation, we can see that our shaping only method achieves higher cumulative reward than the method “shaping without  $r - \hat{r}$ ”. This shows the high fidelity of our high-level explanation.

## F.6 HOW TO ACCELERATE THE TRIPLE-LOOP ALGORITHM

The total iterations of UTILITY is  $N \times \tilde{N} \times \tilde{N}$  where  $N$  is the iteration number of the outer loop,  $\tilde{N}$  is the iteration number of the middle loop, and  $\tilde{N}$  is the iteration number of the inner loop. While the triple-loop structure looks computationally expensive, in practice, we can significantly accelerate the algorithm using warm start in the inner loop and middle loop. Take the inner loop as an example, given the shaping parameter  $\theta$  and the current dual parameter  $\lambda_{\tilde{n}}$ , we need to compute the corresponding constrained soft Bellman policy  $\pi_{\lambda_{\tilde{n}};\theta}$  in the inner loop. Instead of starting from a random policy initialization, we use the policy  $\hat{\pi}_{\lambda_{\tilde{n}-1};\theta}$  learned in last inner loop as the initialization where  $\hat{\pi}_{\lambda_{\tilde{n}-1};\theta}$  is an approximation of  $\pi_{\lambda_{\tilde{n}-1};\theta}$ . The intuition behind this is that since  $\lambda_{\tilde{n}}$  and  $\lambda_{\tilde{n}-1}$  are close (only different by one-step gradient descent), it is expected that  $\pi_{\lambda_{\tilde{n}-1};\theta}$  and  $\pi_{\lambda_{\tilde{n}};\theta}$  are close. Therefore, using  $\hat{\pi}_{\lambda_{\tilde{n}-1};\theta}$  as the initialization makes it easier to approach  $\pi_{\lambda_{\tilde{n}};\theta}$ . Therefore, the warm start trick enables us to use fewer iterations for the inner loop, i.e., the iteration number  $\tilde{N}$  reduces. We use the similar warm start trick for the middle loop to reduce the iteration number  $\tilde{N}$ .

Empirically, in our experiment, we use warm start and set inner iteration number  $\tilde{N} = 1$  and middle iteration number  $\tilde{N} = 2$ . We can see that even if it is a triple-loop algorithm, the total iteration number is small, i.e.,  $2N$ . This warm start trick is inspired by (Zeng et al., 2022; Liu & Zhu, 2024a) where they also use warm start and only run the inner loop for one iteration and the final results are not worse than the algorithm that runs the inner loop for many iterations starting from random initialization. We run the code on a desktop whose CPU is Intel Core i9 12900k and GPU is NVIDIA RTX 3080. We include the runtime of our UTILITY algorithm below:

Table 10: Computation time.

	Delayed HalfCheetah	Delayed Hopper	Delayed Walker	Delayed Ant	HalfCheetah	Hopper	Walker	Ant
UTILITY	2h30min	2h10min	2h32min	3h11min	2h26min	2h5min	2h34min	3h09min
RICE	2h18min	1h53min	2h19min	2h55min	2h18min	1h59min	2h10min	2h58min
SIL	3h14min	3h32min	3h58min	4h43min	3h55min	3h8min	4h11min	4h49min
LIR	2h41min	2h5min	2h29min	3h15min	2h35min	2h13min	2h45min	3h16min

Table 10 shows that the computation time of UTILITY is comparable to the baselines due to the warm start trick.

## G POTENTIAL SOCIETAL IMPACT

This paper has positive impact which is to improve the performance of RL agents. However, the paper also has potential negative impacts. Since the two-level explanation identifies the mistakes made by the RL agents. A malicious entity may use these weaknesses or mistakes to launch attack to the RL agents. To alleviate this issue, one solution is to keep the demonstration of the RL agent private, so that the malicious entity cannot get access to the demonstration and thus cannot find the weakness.

## H LIMITATIONS

One limitation of the method is that it requires to interact with the environment. Therefore, one future work is to extend this method to the offline RL setting.

## I EVOLUTION FRO FIGURE 1C TO FIGURE 1D

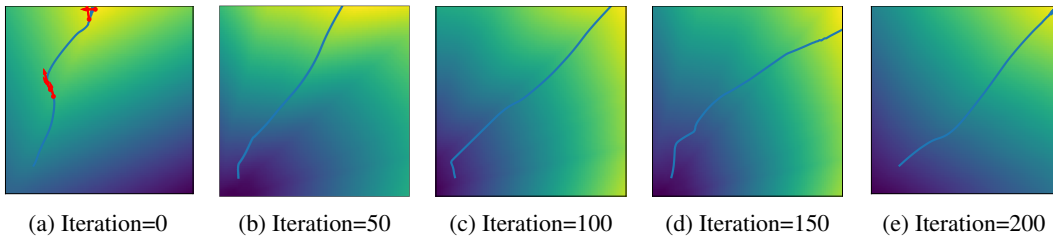


Figure 3: Evolution from Figure 1c to Figure 1d.

For this scenario, we run UTILITY for 200 iterations. Figure 3 show the evolution from Figure 1c (iteration=0) to Figure 1d (iteration=200). From the evolution, we can see that the learned reward becomes closer and closer to the ground truth reward and the learned policy becomes more and more optimal.