
Maximum-Likelihood Inverse Reinforcement Learning with Finite-Time Guarantees

Siliang Zeng¹ Chenliang Li² Alfredo Garcia³ Mingyi Hong¹

Abstract

Inverse reinforcement learning (IRL) aims to recover the reward function and the associated optimal policy that best fits observed sequences of states and actions implemented by an expert. Many algorithms for IRL have an inherent nested structure: the inner loop finds the optimal policy given parametrized rewards while the outer loop updates the estimates towards optimizing a measure of fit. For high dimensional environments such nested-loop structure entails a significant computational burden. To reduce the computational burden of a nested loop, novel methods such as SQIL (Reddy et al., 2020) and IQ-Learn (Garg et al., 2021) emphasize policy estimation at the expense of reward estimation accuracy. However, without accurate estimated rewards, it is not possible to do counterfactual analysis such as predicting the optimal policy under different environment dynamics and/or learning new tasks. In this paper we develop a novel *single-loop* algorithm for IRL that does not compromise reward estimation accuracy. In the proposed algorithm, each policy improvement step is followed by a stochastic gradient step for likelihood maximization. We show that the proposed algorithm provably converges to a stationary solution with a finite-time guarantee. If the reward is parameterized linearly, we show the identified solution corresponds to the solution of the maximum entropy IRL problem. Finally, by using robotics control problems in

Mujoco and their transfer settings, we show that the proposed algorithm achieves superior performance compared with other IRL and imitation learning benchmarks.

1. Introduction

Given observed trajectories of states and actions implemented by an expert, we consider the problem of estimating the reinforcement learning environment in which the expert was trained. This problem is generally referred to as inverse reinforcement learning (IRL) (see (Osa et al., 2018) for a recent survey). Assuming the environment dynamics are known (or available online), the IRL problem consists of estimating the reward function and the expert’s policy (optimizing such rewards) that best fits the data. While there are limitations on the identifiability of rewards (Kim et al., 2021), the estimation of rewards based upon expert trajectories enables important counterfactual analysis such as the estimation of optimal policies under *different* environment dynamics and/or reinforcement learning of *new* tasks.

In the seminal work (Ziebart et al., 2008), the authors developed an IRL formulation, in which the model for the expert’s behavior is the policy that maximizes entropy subject to a constraint requiring that the expected features under such policy match the empirical averages in the expert’s observation dataset. The algorithms developed for MaxEnt-IRL (Ziebart et al., 2008; 2010; Wulfmeier et al., 2015) have a nested loop structure, alternating between an outer loop with a reward update step, and an inner loop that calculates the explicit policy estimates. The computational burden of this nested structure is manageable in tabular environments, but it becomes significant in high dimensional settings requiring function approximation.

Towards developing more efficient IRL algorithms, a number of works (Finn et al., 2016a;b; Ho & Ermon, 2016; Fu et al., 2017) propose to leverage the idea of adversarial training (Goodfellow et al., 2014). These algorithms learn a non-stationary reward function through training a discriminator, which is then used to guide the policy to match the behavior trajectories from the expert dataset. However, (Ni et al., 2020) pointed out that the resulting

¹Department of Electrical and Computer Engineering, University of Minnesota ²School of Data Science, The Chinese University of Hong Kong, Shenzhen ³Department of Industrial and Systems Engineering, Texas A&M University. Correspondence to: Siliang Zeng <zeng0176@umn.edu>, Chenliang Li <chenliangli@link.cuhk.edu.cn>, Alfredo Garcia <alfredo.garcia@tamu.edu>, Mingyi Hong <mhong@umn.edu>.

discriminator (hence the reward function) typically cannot be used in new learning tasks, since it is highly dependent on the corresponding policy and current environment dynamics. Moreover, due to the brittle approximation techniques and sensitive hyperparameter choice in the adversarial training, these IRL algorithms can be unstable. (Kurach et al., 2018; Kostrikov et al., 2019).

More recent works (Reddy et al., 2020; Garg et al., 2021) have developed algorithms to alleviate the computational burden of the nested-loop training procedures. In (Reddy et al., 2020), the authors propose to model the IRL using certain maximum entropy RL problem with specific reward function (which assigns $r = +1$ for matching expert demonstrations and $r = 0$ for all other behaviors). Then a soft Q imitation learning (SQIL) algorithm is developed. In (Garg et al., 2021), the authors propose to transform the standard formulation of IRL (discussed above) into a single-level problem, through learning a soft Q-function to implicitly represent the reward function and the policy. An inverse soft-Q learning (IQ-Learn) algorithm is then developed, which is shown to be effective in estimating the policy for the environment that it is trained on. Despite being computationally efficient, IQ-Learn sacrifices the accuracy in estimating the rewards since it indirectly recovers rewards from a soft Q-function approximator which is highly dependent upon the environment dynamics and does not strictly satisfy the soft-Bellman equation. Therefore it is not well-suited for counterfactual prediction or transfer learning setting.

In f -IRL (Ni et al., 2020) the authors consider an approach for estimating rewards based on the minimization of several measures of divergence with respect to the expert’s state visitation measure. The approach is limited to estimating rewards that only depend on state. Moreover, while the results reported are based upon a single-loop implementation, the paper does not provide a convergence guarantee to support performance.

In the end, we introduce several related works which focus on developing “robust” algorithms and enable the agent to adapt to new environment with different dynamics. There is a line of works focusing on disentangling the reward function from the environment dynamics, so that the recovered reward functions could be transferred across environments with different dynamics. In (Fu et al., 2017), the authors propose an algorithm called adversarial inverse reinforcement learning (AIRL). Constructing the estimated reward as a function which depends on the current state and next state, AIRL enables the agent to learn policies in a new environment through leveraging the estimated reward function recovered from the training environment. In (Kim et al., 2021), the authors prove necessary and sufficient conditions for reward identifiability

in deterministic MDP models with the maximum entropy reinforcement learning objective. In (Cao et al., 2021), the authors present a theoretical analysis to show the necessary and sufficient condition to identify an action-independent time-homogeneous reward function under the maximum entropy IRL problem. A recent line of works consider a more challenging setting, where the learner has no access to the expert environment, and there is a transition dynamics mismatch between the expert and the learner. In (Liu et al., 2020), the authors propose a state alignment based imitation learning method so that the imitator could follow the state sequences in expert demonstrations as much as possible. Arguing that the expert actions are not efficient demonstrations under transition dynamics mismatch, (Gangwani & Peng, 2020) further develops a state-only imitation learning method. In (Viano et al., 2021), the authors revisit the Maximum Causal Entropy IRL when there is a transition dynamics mismatch between the expert and the learner. A theoretical analysis is further provided in (Viano et al., 2021) to show the upper bound on the learner’s performance degradation, which is measured in terms of the ℓ_1 -distance between the transition dynamics of the expert and the learner.

Our Contributions. The goal of this work is to develop an algorithm for IRL which is capable of producing high-quality estimates of *both* rewards and behavior policies with finite-time guarantees. The major contributions of this work are listed below.

- We consider a formulation of IRL based on maximum likelihood (ML) estimation over optimal (entropy-regularized) policies, and prove that a strong duality relationship with maximum entropy IRL holds if rewards are represented by a *linear* combination of features.¹ The ML formulation is a *bi-level* optimization problem, where the upper-level problem maximizes the likelihood function, while the lower-level finds the optimal policy under the current reward parameterization. Such a bi-level structure is not only instrumental to the subsequent algorithm design, but is also flexible to incorporate the use of state-only, as well as the regular reward function (which depends on the state and action pair). The former is suitable for transfer learning since it is *insensitive* to the changes of the environment dynamics, while the latter can be used to efficiently imitate the expert policy.
- Based on the ML-IRL formulation, we develop an efficient algorithm. To avoid the computational burden of repeatedly solving the lower-level policy optimization

¹Heuristic arguments for this duality result are discussed in (Ziebart et al., 2008) wherein the distribution of state-action paths is approximated (see equation (4) in (Ziebart et al., 2008)) and the equivalence between maximum entropy estimation and maximum likelihood (over the class of exponential distributions) (Jaynes, 1957) is invoked.

problem, the proposed algorithm has a single-loop structure where the policy improvement step and reward optimization step are performed alternately so that each step can be performed relatively cheaply. Further, we show that the algorithm has strong theoretical guarantees: to achieve certain ϵ -approximate stationary solution for a non-linearly parameterized problem, it requires $\mathcal{O}(\epsilon^{-2})$ steps of policy and reward updates each. To our knowledge, it is the first algorithm which has finite-time guarantee for the IRL problem under nonlinear parameterization of reward functions.

- We conduct extensive experiments to demonstrate that the proposed algorithm outperforms many state-of-the-art IRL algorithms in *both* policy estimation and reward recovery. In particular, when transferring to a new environment, RL algorithms using rewards recovered by the proposed algorithm outperform those that use rewards recovered from existing IRL and imitation learning benchmarks.

2. Preliminaries

In this section, we review the fundamentals of the maximum entropy inverse reinforcement learning (MaxEnt-IRL). We consider an MDP defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \eta, r, \gamma)$; \mathcal{S} and \mathcal{A} denote the state space and the action space respectively; $\mathcal{P}(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ denotes the transition probability; $\eta(\cdot)$ denotes the distribution for the initial state; $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow R$ is the reward function and γ is a discount factor.

Given a dataset \mathcal{D} which contains trajectories $\tau = \{(s_t, a_t)\}_{t=0}^{\infty}$ sampled from an expert policy, the MaxEnt-IRL formulation (Ziebart et al., 2010; 2013; Bloem & Bambos, 2014) consists of finding a policy maximizing entropy subject to the expected features under such policy matching the empirical averages in the expert's observation dataset. Specifically, the MaxEnt-IRL formulation is given by:

$$\begin{aligned} \max_{\pi} H(\pi) &:= \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \right] \\ &\quad \text{(MaxEnt-IRL)} \\ \text{s.t. } \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] &= \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \end{aligned} \quad (1)$$

where $\phi(s, a)$ is the feature vector of the state-action pair (s, a) . Let θ denote the dual variable for the linear constraint, then the Lagrangian of (MaxEnt-IRL) is given by

$$\begin{aligned} \mathcal{L}(\pi, \theta) &:= H(\pi) + \theta, \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \theta^T \phi(s_t, a_t) \right] \\ &\quad - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \theta^T \phi(s_t, a_t) \right]. \end{aligned} \quad (2)$$

In (Ziebart et al., 2010; 2013; Bloem & Bambos, 2014), the authors proposed a "dual descent" algorithm, which alternates between *i*) solving $\max_{\pi} \mathcal{L}(\pi, \theta)$ for fixed θ , and *ii*) a gradient descent step to optimize the dual variable θ . It is shown that the optimizer π_{θ}^* in step *i*) can be recursively defined as $\pi_{\theta}^*(a_t | s_t) = \frac{Z_{a_t | s_t, \theta}}{Z_{s_t, \theta}}$, where $\log Z_{a_t | s_t, \theta} = \phi(s_t, a_t)^T \theta + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)} [\log Z_{s_{t+1}, \theta}]$ and $\log Z_{s_t, \theta} = \log \left(\sum_{a \in \mathcal{A}} Z_{a | s_t, \theta} \right)$.

From a computational perspective, the above algorithm is not efficient: it has a nested-loop structure, which repeatedly computes the optimal policy π_{θ}^* under each variable θ . It is known that when the underlying MDP is of high-dimension, such an algorithm can be computationally prohibitive (Finn et al., 2016b; Ho & Ermon, 2016).

Recent work (Garg et al., 2021) proposed an algorithm called IQ-Learn to improve upon the MaxEnt-IRL by considering a saddle-point formulation:

$$\begin{aligned} \min_r \max_{\pi} \left\{ H(\pi) + \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \cdot r(s_t, a_t) \right] \right. \\ \left. - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \cdot r(s_t, a_t) \right] \right\}. \end{aligned} \quad (3)$$

The authors show that this problem can be transformed into an optimization problem *only* defined in terms of the soft Q-function, which implicitly represents both reward and policy. IQ-Learn is shown to be effective in imitating the expert behavior while only relying on the estimation of the soft Q-function. However, the implicit reward estimate obtained is not necessarily accurate since its soft Q-function estimate depends on the environment dynamics and does not strictly satisfy the soft-Bellman equation. Hence, it is difficult to transfer the recovered reward function to new environments.

3. Problem Formulation

In this section, we consider a ML formulation of the IRL problem and formalize a duality relationship with maximum entropy-based formulation (MaxEnt-IRL).

Maximum Log-Likelihood IRL (ML-IRL)

Let $\mathcal{D} := \{\tau_i\}_{i=1}^N$ denote the dataset containing observed trajectories. A model of the expert's behavior is a randomized policy $\pi_{\theta}(\cdot | s)$ where θ is a parameter vector. Assuming the state dynamics $\mathcal{P}(s_{t+1} | s_t, a_t)$ are known, the discounted log-likelihood of observing a sample trajectory τ under model π_{θ} can be written follows:

$$\begin{aligned} \mathbb{E}_{\tau \sim \mathcal{D}} \left[\log \prod_{t \geq 0} (\mathcal{P}(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t))^{\gamma^t} \right] \\ = \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t (\log \pi_{\theta}(a_t | s_t) + \log \mathcal{P}(s_{t+1} | s_t, a_t)) \right]. \end{aligned}$$

Let $L(\theta) := \mathbb{E}_{\tau \sim \mathcal{D}} [\sum_{t \geq 0} \gamma^t \log \pi_\theta(a_t | s_t)]$. We consider the following maximum log-likelihood IRL formulation:

$$\begin{aligned} & \max_{\theta} L(\theta) && \text{(ML-IRL)} \\ \text{s.t. } & \pi_\theta := \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t; \theta) + \mathcal{H}(\pi(\cdot | s_t)) \right) \right] \end{aligned} \quad (4a)$$

where $r(s, a; \theta)$ is the reward function and $\mathcal{H}(\pi(\cdot | s)) := -\sum_{a \in \mathcal{A}} \pi(a | s) \log \pi(a | s)$ denotes the entropy of policy $\pi(\cdot | s)$.

We now make some remarks about ML-IRL. First, the problem takes the form of a *bi-level* optimization problem, where the *upper-level* problem (ML-IRL) optimizes the reward parameter θ , while the *lower-level* problem describes the expert’s policy as the solution to an entropy-regularized MDP ((Haarnoja et al., 2017; 2018)). In what follows we will leverage recently developed (stochastic) algorithms for bi-level optimization (Hong et al., 2020; Ji et al., 2021; Khanduri et al., 2021), that avoid the high complexity resulted from nested loop algorithms. Second, it is reasonable to use the ML function as the loss, because it searches for a reward function which generates a behavior policy that can best fit the expert demonstrations. While the ML function has been considered in (Jain et al., 2019; Sanghvi et al., 2021), they rely on heuristic algorithms with nested-loop computations to solve their IRL formulations, and the theoretical properties are not studied. Finally, the lower-level problem has been well-studied in the literature (Haarnoja et al., 2017; 2018; Cen et al., 2021; Cayci et al., 2021; Nachum et al., 2017). The entropy regularization in (4a) ensures the uniqueness of the optimal policy π_θ under the fixed reward function $r(s, a; \theta)$ (Haarnoja et al., 2017; Cen et al., 2021). Even when the underlying MDP is high-dimensional and/or complex, the optimal policy could still be obtained; see recent developments in (Haarnoja et al., 2017; 2018). Subsequently, we will leverage these algorithms, together with recent advances in bi-level optimization to develop solution methods for (ML-IRL).

We close this section by formally establishing a connection between (MaxEnt-IRL) and (ML-IRL).

Theorem 3.1. (Strong Duality) *Suppose that the reward function is given as: $r(s, a; \theta) := \phi(s, a)^T \theta$, for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$. Then (ML-IRL) is the Lagrangian dual of (MaxEnt-IRL). Furthermore, strong duality holds, that is: $L(\theta^*) = H(\pi^*)$, where θ^* and π^* are the global optimal solutions for problems (ML-IRL) and (MaxEnt-IRL), respectively.*

The proof of Theorem 3.1 is relegated to Appendix F. To our knowledge this result which specifically addresses the (MaxEnt-IRL) formulation is novel. A duality between

ML estimation and maximum causal entropy is obtained in (Ziebart et al., 2013, Theorem 3). However, the problem considered in that paper is not in RL nor IRL setting, therefore they cannot be directly used in the context of the present paper.

The above duality result reveals a strong connection between the two formulations under linear reward parameterization. Of course, when non-linear parameterization is used strong duality may not hold. However, this connection still suggests that ML-IRL formulation is a meaningful framework for IRL.

4. The Proposed Algorithm

In this section, we design algorithms for (ML-IRL). Recall that one major drawback of algorithms for (MaxEnt-IRL) is that, they repeatedly solve certain policy optimization problem in the inner loop. For example, the algorithm IQ-Learn (Garg et al., 2021) improves the computational efficiency through implicitly representing the reward function and the policy by a Q-function approximator, while the estimation accuracy of the recovered reward is sacrificed. Therefore, one important goal of our design is to find provably efficient algorithms that can avoid high-complexity operations and accurately recover the reward function. Specifically, it is desirable that the resulting algorithm only uses a finite number of reward and policy updates to reach certain high-quality solutions.

To proceed, we will leverage the special *bi-level* structure of the ML-IRL problem. The idea is to alternate between one step of policy update to improve the solution of the lower-level problem, and one step of the parameter update which improves the upper-level loss function. At each iteration k , given the current policy π_k and the reward parameter θ_k , a new policy π_{k+1} is generated from the policy improvement step, and θ_{k+1} is generated by the reward optimization step.

This kind of alternating update is efficient, because there is no need to completely solve the policy optimization subproblem, before updating the reward parameters. It has been used in many other RL related settings as well. For example, the well-known actor-critic (AC) algorithm for policy optimization (Konda & Tsitsiklis, 1999; Wu et al., 2020; Hong et al., 2020) alternates between one step of policy update, and one step of critic parameter update. Below we present the details of our algorithm at a given iteration k .

Policy Improvement Step. Given the reward parameter θ_k is fixed, let us consider optimizing the lower-level problem. Towards this end, define the so-called soft Q and soft value

Algorithm 1 Maximum Likelihood Inverse Reinforcement Learning (ML-IRL)

- 0: **Input:** Initialize reward parameter θ_0 and policy π_0 . Set the reward parameter's stepsize as α .
- 0: **for** $k = 0, 1, \dots, K - 1$ **do**
- 0: **Policy Evaluation:** Compute $Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(\cdot, \cdot)$ under reward function $r(\cdot, \cdot; \theta_k)$
- 0: **Policy Improvement:** $\pi_{k+1}(\cdot|s) \propto \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \cdot)), \forall s \in \mathcal{S}$.
- 0: **Data Sampling I:** Sampling expert trajectory $\tau_k^E := \{s_t, a_t\}_{t \geq 0}$ from the dataset D
- 0: **Data Sampling II:** Sampling agent trajectory $\tau_k^A := \{s_t, a_t\}_{t \geq 0}$ from the policy π_{k+1}
- 0: **Estimating Gradient:** $g_k := h(\theta_k; \tau_k^E) - h(\theta_k; \tau_k^A)$ where $h(\theta; \tau) := \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta)$
- 0: **Reward Parameter Update:** $\theta_{k+1} := \theta_k + \alpha g_k$

functions for a given policy π_k and a reward parameter θ_k :

$$V_{r_k, \pi_k}^{\text{soft}}(s) = \mathbb{E}_{\pi_k, s_0=s} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t; \theta_k) + \mathcal{H}(\pi_k(\cdot|s_t)) \right) \right] \quad (5a)$$

$$Q_{r_k, \pi_k}^{\text{soft}}(s, a) = r(s, a; \theta_k) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} [V_{r_k, \pi_k}^{\text{soft}}(s')] \quad (5b)$$

We will adopt the well-known *soft policy iteration* (Haarnoja et al., 2017) to optimize the lower-level problem (4a). Under the current reward parameter θ_k and the policy π_k , the soft policy iteration generates a new policy π_{k+1} as follows

$$\pi_{k+1}(a|s) \propto \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a)), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}. \quad (6)$$

Under a fixed reward function, it can be shown that the new policy π_{k+1} monotonically improves π_k , and it converges linearly to the optimal policy; see (Haarnoja et al., 2017, Theorem 4) and (Cen et al., 2021, Theorem 1).

Note that in practice, we usually do not have direct access to the exact soft Q-function in (5b). In order to perform the policy improvement, a few stochastic update steps in soft Q-learning (Haarnoja et al., 2017) or soft Actor-Critic (SAC) (Haarnoja et al., 2018) could be used to replace the one-step soft policy iteration (6). In the appendix, we present Alg. 2 to demonstrate such alternative implementation of our proposed algorithm.

Reward Optimization Step. We propose to use a stochastic gradient-type algorithm to optimize θ . Towards this end, let us first derive the exact gradient $\nabla L(\theta)$. See Appendix C for detailed proof.

Lemma 4.1. *The gradient of the likelihood function $L(\theta)$ can be expressed as follows:*

$$\begin{aligned} \nabla L(\theta) &= \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right] \\ &\quad - \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right]. \end{aligned} \quad (7)$$

To obtain stochastic estimators of the exact gradient $\nabla L(\theta_k)$, we take two approximation steps: 1) approximate

the optimal policy π_{θ_k} by π_{k+1} in (6), since the optimal policy π_{θ_k} is not available throughout the algorithm; 2) sample one expert trajectory τ_k^E from the dataset \mathcal{D} and one agent trajectory τ_k^A from the current policy π_{k+1} so to approximate the expectation operators in (7).

Following the approximation steps mentioned above, we construct a stochastic estimator g_k to approximate the exact gradient $\nabla L(\theta_k)$ in (7) as follows:

$$g_k = h(\theta_k; \tau_k^E) - h(\theta_k; \tau_k^A) \quad (8)$$

where $h(\theta; \tau) = \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta)$. With the stochastic gradient estimator g_k , the reward parameter θ_k is updated as:

$$\theta_{k+1} = \theta_k + \alpha g_k. \quad (9)$$

where α is the stepsize in updating the reward parameter.

In summary, the proposed algorithm for solving the ML-IRL problem (ML-IRL) is given in Alg. 1.

5. Theoretical Analysis

In this section, we present finite-time guarantees for the proposed algorithm.

To begin with, first recall that in Sec. 3, we have mentioned that (ML-IRL) is a bi-level problem, where the upper level (resp. the lower level) problem optimizes the reward parameter (resp. the policy). In order to solve (ML-IRL), our algorithm 1 has a *single-loop* structure, which alternates between one step of policy update and one step of the reward parameter update. Such a single-loop structure indeed has computational benefit, but it also leads to potential unstableness, since the lower level problem can stay far away from its true solutions. Specifically, at each iteration k , the potential unstableness is induced by the distribution mismatch between the policy π_{k+1} and π_{θ_k} , when we use estimator g_k (8) to approximate the exact gradient $\nabla L(\theta_k)$ (7) in updating the reward parameter θ_k .

Towards stabilizing the algorithm, we adopt the so-called *two-timescale* stochastic approximation (TTSA) approach (Borkar, 1997; Hong et al., 2020), where the lower-level

problem updates in a faster time-scale (i.e., converges faster) compared with its upper-level counterpart. Intuitively, the TTSA enables the π_{k+1} tracks the optimal π_{θ_k} , leading to a stable algorithm. In the proposed Algorithm 1, the policy (lower-level variable) is continuously updated by the soft policy iteration (6), and it is ‘fast’ because it converges linearly to the optimal policy under a fixed reward function (Theorem 1]cen2021fast. On the other hand, the reward parameter update (9) does not have such linear convergence property, therefore it works in a ‘slow’ timescale. To begin our analysis, let us first present a few technical assumptions.

Assumption 5.1 (Ergodicity). For any policy π , assume the Markov chain with transition kernel \mathcal{P} is irreducible and aperiodic under policy π . Then there exist constants $\kappa > 0$ and $\rho \in (0, 1)$ such that

$$\sup_{s \in \mathcal{S}} \|\mathbb{P}(s_t \in \cdot | s_0 = s, \pi) - \mu_\pi(\cdot)\|_{TV} \leq \kappa \rho^t, \quad \forall t \geq 0$$

where $\|\cdot\|_{TV}$ is the total variation (TV) norm; μ_π is the stationary state distribution under π .

Assumption 5.1 assumes the Markov chain mixes at a geometric rate. It is a common assumption in the literature of RL (Bhandari et al., 2018; Zou et al., 2019; Wu et al., 2020), which holds for any time-homogeneous Markov chain with finite-state space or any uniformly ergodic Markov chain with general state space.

Assumption 5.2. For any $s \in \mathcal{S}$, $a \in \mathcal{A}$ and any reward parameter θ , the following holds:

$$\|\nabla_{\theta} r(s, a; \theta)\| \leq L_r, \quad (10a)$$

$$\|\nabla_{\theta} r(s, a; \theta_1) - \nabla_{\theta} r(s, a; \theta_2)\| \leq L_g \|\theta_1 - \theta_2\| \quad (10b)$$

where L_r and L_g are positive constants.

Assumption 5.2 assumes that the parameterized reward function has bounded gradient and is Lipschitz smooth. Such assumption in Lipschitz property are common in the literature of min-max / bi-level optimization (Jin et al., 2020; Hong et al., 2020; Chen et al., 2021; Guan et al., 2021; Khanduri et al., 2021).

Based on Assumptions 5.1 - 5.2, we next provide the following Lipschitz properties:

Lemma 5.3. *Suppose Assumptions 5.1 - 5.2 hold. Given any reward parameters θ_1 and θ_2 , the following results hold for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$:*

$$|Q_{r_{\theta_1}, \pi_{\theta_1}}^{\text{soft}}(s, a) - Q_{r_{\theta_2}, \pi_{\theta_2}}^{\text{soft}}(s, a)| \leq L_q \|\theta_1 - \theta_2\|, \quad (11a)$$

$$\|\nabla L(\theta_1) - \nabla L(\theta_2)\| \leq L_c \|\theta_1 - \theta_2\| \quad (11b)$$

where $Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(\cdot, \cdot)$ denotes the soft Q-function under the reward function $r(\cdot, \cdot; \theta)$ and the policy π_{θ} . The positive constants L_q and L_c are defined in Appendix D.

The Lipschitz properties identified in Lemma 5.3 are vital for the convergence analysis. Then we present the main results, which show the convergence speed of the policy $\{\pi_k\}_{k \geq 0}$ and the reward parameter $\{\theta_k\}_{k \geq 0}$ in the Alg. 1. Please see Appendix D for the detailed proof.

Theorem 5.4. *Suppose Assumptions 5.1 - 5.2 hold. Selecting stepsize $\alpha := \frac{\alpha_0}{K^\sigma}$ for the reward update step (9) where $\alpha_0 > 0$ and $\sigma \in (0, 1)$ are some fixed constants, and K is the total number of iterations to be run by the algorithm. Then the following result holds:*

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty}] = \mathcal{O}(K^{-1}) + \mathcal{O}(K^{-\sigma}) \quad (12a)$$

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla L(\theta_k)\|^2] = \mathcal{O}(K^{-\sigma}) + \mathcal{O}(K^{-1+\sigma}) + \mathcal{O}(K^{-1}) \quad (12b)$$

where we denote $\|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty} := \max_{s \in \mathcal{S}, a \in \mathcal{A}} |\log \pi_{k+1}(a|s) - \log \pi_{\theta_k}(a|s)|$. In particular, setting $\sigma = 1/2$, then both quantities in (12a) and (12b) converge with the rate $\mathcal{O}(K^{-1/2})$.

In Theorem 5.4, we present the finite-time guarantee for the convergence of the Alg.1. As a remark, we note that our theoretical guarantee is different from the existing results, such as (Cen et al., 2021) that showed the convergence rate of soft policy iteration under a fixed reward function. Theorem 5.4 analyzes a more challenging setting where both the policy and reward parameter are kept changing. To present the key steps in our analysis, we provide a proof sketch below. The detailed proof is in Appendix G.

Proof sketch. We outline our main steps in analyzing (12a) and (12b) respectively.

In order to show the convergence of policy estimates in (12a), there are several key steps. First, we note that both policies π_{k+1} and π_{θ_k} are in the softmax parameterization, where $\pi_{k+1}(\cdot|s) \propto \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \cdot))$ and $\pi_{\theta_k}(\cdot|s) \propto \exp(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, \cdot))$. Then, we can show a Lipschitz continuity property between the policy and the soft Q-function:

$$\|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty} \leq 2 \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty},$$

where the infinity norm $\|\cdot\|_{\infty}$ is defined over the state-action space $\mathcal{S} \times \mathcal{A}$. Moreover, by analyzing the contraction property of the soft policy iteration (6), we bound $\|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty}$ as:

$$\begin{aligned} & \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \\ & \leq \gamma \|Q_{r_{\theta_{k-1}}, \pi_{k-1}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + 2L_q \|\theta_k - \theta_{k-1}\|. \end{aligned}$$

To ensure that the error term $\|\theta_k - \theta_{k-1}\|$ is small, we select the stepsize of reward parameters as $\alpha := \frac{\alpha_0}{K\sigma}$, where K is the total number of iterations and $\sigma > 0$. Then, by combining previous two steps, we could further show the convergence rate of the policy estimates in (12a).

To prove the convergence of the reward parameters in (12b), we first leverage the Lipschitz smooth property of $L(\theta)$ in (11b). However, one technical challenge in the convergence analysis is how to handle the bias between the gradient estimator g_k defined in (8) and the exact gradient $\nabla L(\theta_k)$. When we construct the gradient estimator g_k in (8), we need to sample trajectories from the current policy π_{k+1} and the expert dataset \mathcal{D} . However, according to the expression of $\nabla L(\theta_k)$ in (7), the trajectories are sampled from the optimal policy π_{θ_k} and the dataset \mathcal{D} . Hence, there is a distribution mismatch between π_{k+1} and π_{θ_k} . Our key idea is to leverage (12a) to handle this distribution mismatch error, and thus show that the bias between g_k and $\nabla L(\theta_k)$ could be controlled. \square

6. A Discussion over State-Only Reward

In this section we consider the IRL problems modeled by using rewards that are only a function of the state. A lower dimensional representation of the agent’s preferences (i.e. in terms only of states as opposed to states *and* actions) is more likely to facilitate counterfactual analysis such as predicting the optimal policy under different environment dynamics and/or learning new tasks. This is because the estimation of preferences which are only defined in terms of states is less sensitive to the specific environment dynamics in the expert’s demonstration dataset. Moreover, in application such as healthcare (Yu et al., 2021) and autonomous driving (Kiran et al., 2021), where simply imitating the expert policy can potentially result in poor performance, since the learner and the expert may have different transition dynamics. Similar points have also been argued in recent works (Gangwani & Peng, 2020; Ni et al., 2020; Viano et al., 2021).

Next, let us briefly discuss how we can understand (ML-IRL) and Alg.1, when the reward is parameterized as a state-only function.

Lemma 6.1. *Suppose the expert trajectories τ is sampled from a policy π^E , and the reward is parameterized as a state-only function $r(s; \theta)$. Then ML-IRL is equivalent to the following:*

$$\begin{aligned} \min_{\theta} \mathbb{E}_{s_0 \sim \eta(\cdot)} [V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0) - V_{r_{\theta}, \pi^E}^{\text{soft}}(s_0)] \quad (13a) \\ \text{s.t. } \pi_{\theta} := \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t; \theta) + \mathcal{H}(\pi(\cdot|s_t)) \right) \right]. \end{aligned}$$

Please see Appendix E for the detailed derivation.

Intuitively, the above lemma says that, when dealing with the state-only IRL, (ML-IRL) minimizes the gap between the soft value functions of the optimal policy π_{θ} and that of the expert policy π^E . Moreover, Alg.1 can also be easily implemented with the state-only reward. In fact, the entire algorithm essentially stays the same, and the only change is that $r(s, a; \theta)$ will be replaced by $r(s; \theta)$. Therefore, even under the state-only IRL setting where the expert dataset only contains visited states, our formulation and the proposed algorithm still work if we parameterize the reward as a state-only function.

7. Numerical Results

In this section, we test the performance of our algorithm on a diverse collection of RL tasks and environments. In each experiment set, we train algorithms until convergence and average the scores of the trajectories over multiple random seeds. The hyperparameter settings and simulation details are provided in Appendix A.

Mujoco Tasks For Inverse Reinforcement Learning.

In this experiment set, we test the performance of our algorithm on imitating the expert behavior. We consider several high-dimensional robotics control tasks in Mujoco (Todorov et al., 2012). Two class of existing algorithms are considered as the comparison baselines: 1) imitation learning algorithms that only learn the policy to imitate the expert, including Behavior Cloning (BC) (Pomerleau, 1988) and Generative Adversarial Imitation Learning (GAIL) (Ho & Ermon, 2016); 2) IRL algorithms which learn a reward function *and* a policy simultaneously, including Adversarial Inverse Reinforcement Learning (AIRL) (Fu et al., 2017), f -IRL (Ni et al., 2020) and IQ-Learn (Garg et al., 2021). To ensure fair comparison, all imitation learning / IRL algorithms use soft Actor-Critic (Haarnoja et al., 2018) as the base RL algorithm. For the expert dataset, we use the data provided in the official implementation² of f -IRL.

In this experiment, we implement two versions of our proposed algorithm: ML-IRL(State-Action) where the reward is parameterized as a function of state and action; ML-IRL(State-Only) which utilizes the state-only reward function. In Table 1, we present the simulation results under a limited data regime where the expert dataset only contains a single expert trajectory. The scores (cumulative rewards) reported in the table is averaged over 3 random seeds. In each random seed, we train algorithm from initialization and collect 20 trajectories to average their cumulative rewards after the algorithms converge. According to the results reported in Table 1, it shows that our proposed algorithms outperform the baselines on most tasks.

We observe that BC fails to imitate the expert’s behavior. It

²<https://github.com/twni2016/f-IRL>

Maximum-Likelihood Inverse Reinforcement Learning with Finite-Time Guarantees

Task	BC	GAIL	IQ-Learn	f -IRL	ML-IRL (State-Only)	ML-IRL (State-Action)	Expert
Hopper	20.76	2820.44	3299.01	3361.02	3382.44	3392.51	3592.63
Half-Cheetah	-2.05	4812.55	5043.26	4802.33	5254.84	4519.74	5098.3
Walker	-13.84	1183.94	5088.42	4798.01	5373.27	5393.89	5344.21
Ant	990.71	3962.71	4362.90*	4960.18	5074.35	5546.71	5926.18
Humanoid	75.82	4194.38	5227.10*	5441.42	5496.87	5420.18	5351.08

Table 1. **Mujoco Results.** The performance of benchmark algorithms under a single expert trajectory.

Setting	IQ-Learn	AIRL	f -IRL	ML-IRL(State-Only)	Groud-Truth
Data Transfer	-11.78	-5.39	188.85	221.51	320.15
Reward Transfer	-1.04	130.3	156.45	187.69	320.15

Table 2. **Transfer Learning.** The performance of benchmark algorithms under a single expert trajectory. The scores in the table are obtained similarly as in Table 1.

is due to the fact that BC is based on supervised learning and thus could not learn a good policy under such a limited data regime. Moreover, we notice the training of IQ-Learn is unstable, which may be due to its inaccurate approximation to the soft Q-function. Therefore, in the Mujoco tasks where IQ-Learn does not perform well, so that we cannot match the results presented in the original paper (Garg et al., 2021), we directly report results from there (and mark them by * in Table 1). The results of AIRL are not presented in Table 1 since it performs poorly even after spending significant efforts in parameter tuning (similar observations have been made in in (Liu et al., 2020; Ni et al., 2020)).

Transfer Learning Across Changing Dynamics. We further evaluate IRL algorithms on the transfer learning setting. We follows the environment setup in (Fu et al., 2017), where two environments with different dynamics are considered: Custom-Ant vs Disabled-Ant. We compare ML-IRL(State-Only) with several existing IRL methods: 1) AIRL (Fu et al., 2017), 2) f -IRL (Ni et al., 2020); 3) IQ-Learn (Garg et al., 2021).

We consider two transfer learning settings: 1) data transfer; 2) reward transfer. For both settings, the expert dataset / trajectories are generated in Custom-Ant. In the data transfer setting, we train IRL agents in Disabled-Ant by using the expert trajectories, which are generated in Custom-Ant. In the reward transfer setting, we first use IRL algorithms to infer the reward functions in Custom-Ant, and then transfer these recovered reward functions to Disabled-Ant for further evaluation. In both settings, we also train SAC with the ground-truth reward in Disabled-Ant and report the scores.

The numerical results are reoprted in Table 2. the proposed ML-IRL(State-Only) achieves superior performance compared with the existing IRL benchmarks in both settings. We notice that IQ-Learn fails in both settings

since it indirectly recovers the reward function from a soft Q-function approximator, which could be inaccurate and is highly dependent upon the environment dynamics. Therefore, the reward function recovered by IQ-Learn can not be disentangled from the expert actions and environment dynamics, which leads to its failures in the transfer learning tasks.

8. Conclusion

In this paper, we present a maximum likelihood IRL formulation and propose a provably efficient algorithm with a single-loop structure. To our knowledge, we provide the first non-asymptotic analysis for IRL algorithm under nonlinear reward parameterization. As a by-product, when we parameterize the reward as a state-only function, our algorithm could work in state-only IRL setting and enable reward transfer to new environments with different dynamics. Our algorithm outperforms existing IRL methods on high-dimensional robotics control tasks and corresponding transfer learning settings. A limitation of our method is the requirement for online training, so one future direction of this work is to further extend our algorithm and the theoretical analysis to the offline IRL setting.

Potential Negative Social Impacts. Since IRL methods aim to recover the reward function and the associated optimal policy from the observed expert dataset, potential negative social impacts may occur if there are bad demonstrations included in the expert dataset. Thus, for sensitive applications such as autonomous driving and clinical decision support, additional care should be taken to avoid negative biases from the expert demonstrations and ensure safe adaptation.

References

- Bhandari, J., Russo, D., and Singal, R. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on learning theory*, pp. 1691–1692. PMLR, 2018.
- Bloem, M. and Bambos, N. Infinite time horizon maximum causal entropy inverse reinforcement learning. In *53rd IEEE conference on decision and control*, pp. 4911–4916. IEEE, 2014.
- Borkar, V. S. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5):291–294, 1997.
- Cao, H., Cohen, S., and Szpruch, L. Identifiability in inverse reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Cayci, S., He, N., and Srikant, R. Linear convergence of entropy-regularized natural policy gradient with linear function approximation. *arXiv preprint arXiv:2106.04096*, 2021.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 2021.
- Chen, T., Sun, Y., and Yin, W. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing Systems*, 34:25294–25307, 2021.
- Finn, C., Christiano, P., Abbeel, P., and Levine, S. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint arXiv:1611.03852*, 2016a.
- Finn, C., Levine, S., and Abbeel, P. Guided cost learning: Deep inverse optimal control via policy optimization. In *International Conference on Machine Learning*, pp. 49–58. PMLR, 2016b.
- Fu, J., Luo, K., and Levine, S. Learning robust rewards with adversarial inverse reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- Gangwani, T. and Peng, J. State-only imitation with transition dynamics mismatch. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJgLLyrYwB>.
- Garg, D., Chakraborty, S., Cundy, C., Song, J., and Ermon, S. Iq-learn: Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34, 2021.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2014.
- Guan, Z., Xu, T., and Liang, Y. When will generative adversarial imitation learning algorithms attain global convergence. In *International Conference on Artificial Intelligence and Statistics*, pp. 1117–1125. PMLR, 2021.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pp. 1352–1361. PMLR, 2017.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pp. 1861–1870. PMLR, 2018.
- Ho, J. and Ermon, S. Generative adversarial imitation learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- Hong, M., Wai, H.-T., Wang, Z., and Yang, Z. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020.
- Jain, V., Doshi, P., and Banerjee, B. Model-free irl using maximum likelihood estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3951–3958, 2019.
- Jaynes, E. T. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- Ji, K., Yang, J., and Liang, Y. Bilevel optimization: Convergence analysis and enhanced design. In *International Conference on Machine Learning*, pp. 4882–4892. PMLR, 2021.
- Jin, C., Netrapalli, P., and Jordan, M. What is local optimality in nonconvex-nonconcave minimax optimization? In *International Conference on Machine Learning*, pp. 4880–4889. PMLR, 2020.
- Khanduri, P., Zeng, S., Hong, M., Wai, H.-T., Wang, Z., and Yang, Z. A near-optimal algorithm for stochastic bilevel optimization via double-momentum. *Advances in Neural Information Processing Systems*, 34, 2021.
- Kim, K., Garg, S., Shiragur, K., and Ermon, S. Reward identification in inverse reinforcement learning. In *International Conference on Machine Learning*, pp. 5496–5505. PMLR, 2021.

- Kiran, B. R., Sobh, I., Talpaert, V., Mannion, P., Al Sallab, A. A., Yogamani, S., and Pérez, P. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- Konda, V. and Tsitsiklis, J. Actor-critic algorithms. *Advances in Neural Information Processing Systems*, 12, 1999.
- Kostrikov, I., Agrawal, K. K., Dwibedi, D., Levine, S., and Tompson, J. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Hk4fpoA5Km>.
- Kurach, K., Lucic, M., Zhai, X., Michalski, M., and Gelly, S. The gan landscape: Losses, architectures, regularization, and normalization. 2018.
- Liu, F., Ling, Z., Mu, T., and Su, H. State alignment-based imitation learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rylrdxHFDr>.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. Bridging the gap between value and policy based reinforcement learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- Ni, T., Sikchi, H., Wang, Y., Gupta, T., Lee, L., and Eysenbach, B. f-irl: Inverse reinforcement learning via state marginal matching. *arXiv preprint arXiv:2011.04709*, 2020.
- Osa, T., Pajarinen, J., Neumann, G., Bagnell, J. A., Abbeel, P., and Peters, J. *An Algorithmic Perspective on Imitation Learning*, volume 7 of *Foundations and Trends in Robotics*. 2018.
- Pomerleau, D. A. ALVINN: An autonomous land vehicle in a neural network. *Advances in Neural Information Processing Systems*, 1, 1988.
- Reddy, S., Dragan, A. D., and Levine, S. SQIL: Imitation learning via reinforcement learning with sparse rewards. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S1xKd24twB>.
- Sanghvi, N., Usami, S., Sharma, M., Groeger, J., and Kitani, K. Inverse reinforcement learning with explicit policy estimates. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 9472–9480, 2021.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033. IEEE, 2012.
- Viano, L., Huang, Y.-T., Kamalaruban, P., Weller, A., and Cevher, V. Robust inverse reinforcement learning under transition dynamics mismatch. *Advances in Neural Information Processing Systems*, 34, 2021.
- Wu, Y. F., Zhang, W., Xu, P., and Gu, Q. A finite-time analysis of two time-scale actor-critic methods. *Advances in Neural Information Processing Systems*, 33: 17617–17628, 2020.
- Wulfmeier, M., Ondruska, P., and Posner, I. Maximum entropy deep inverse reinforcement learning. *arXiv preprint arXiv:1507.04888*, 2015.
- Xu, T., Wang, Z., and Liang, Y. Improving sample complexity bounds for (natural) actor-critic algorithms. *Advances in Neural Information Processing Systems*, 33: 4358–4369, 2020.
- Yu, C., Liu, J., Nemati, S., and Yin, G. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- Ziebart, B. D., Maas, A. L., Bagnell, J. A., Dey, A. K., et al. Maximum entropy inverse reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.
- Ziebart, B. D., Bagnell, J. A., and Dey, A. K. Modeling interaction via the principle of maximum causal entropy. In *International Conference on Machine Learning*, 2010.
- Ziebart, B. D., Bagnell, J. A., and Dey, A. K. The principle of maximum causal entropy for estimating interacting processes. *IEEE Transactions on Information Theory*, 59(4):1966–1980, 2013.
- Zou, S., Xu, T., and Liang, Y. Finite-sample analysis for sarsa with linear function approximation. *Advances in Neural Information Processing Systems*, 32, 2019.

Algorithm 2 *Practical Implementation of ML-IRL*

```

0: Input: Initialize reward parameter  $\theta_0$  and policy  $\pi_0$ . Set the reward parameter's stepsize as  $\alpha$ .
0: for  $k = 0, 1, \dots, K - 1$  do
0:   Policy Update:  $\pi_{k+1} \leftarrow$  several SAC steps under reward function  $r(\cdot, \cdot; \theta_k)$  and policy  $\pi_k$ .
0:   Data Sampling I: Sampling expert trajectory  $\tau_k^E := \{s_t, a_t\}_{t \geq 0}$  from the dataset  $D$ 
0:   Data Sampling II: Sampling agent trajectory  $\tau_k^A := \{s_t, a_t\}_{t \geq 0}$  from the policy  $\pi_{k+1}$ 
0:   Estimating Gradient:  $g_k := h(\theta_k; \tau_k^E) - h(\theta_k; \tau_k^A)$  where  $h(\theta; \tau) := \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta)$ 
0:   Reward Parameter Update:  $\theta_{k+1} := \theta_k + \alpha g_k$ 
0: end for=0

```

Appendix

A. Experiment Details

A.1. Mujoco Tasks For Inverse Reinforcement Learning.

In all experiments, we test the performance of benchmark algorithms on Hopper, Half-Cheetah, Walker, Ant, Humanoid environments from OpenAI Gym. To ensure fair comparison, we use an open-source implementation³ of SAC as the base RL algorithm for all imitation learning / IRL methods. Moreover, Adam is used as the optimizer in SAC.

In SAC, both policy network and Q-network are (64, 64) MLPs with ReLU activation function, and we set their stepsizes as 3×10^{-3} . Moreover, in our proposed algorithms, we parameterize the reward function by a (64, 64) MLPs with ReLU activation function. For the reward network, we use Adam as the optimizer and the stepsize is set to be 1×10^{-4} .

We present the practical implementation procedure of our proposed algorithm in Table.2. At each iteration, we first warm-start both policy network and Q-network in SAC by using the trained neural networks from the previous iteration. Then, we run 10 episodes in the corresponding Mujoco environment to train the policy network and Q-network in SAC. After that, we collect agent trajectory and expert trajectory to train the reward network by one gradient update.

For the imitation learning / IRL benchmark algorithms, we use their open-source implementations in our experiments. The official implementations of f -IRL is provided in <https://github.com/twni2016/f-IRL>. The official code base for IQ-Learn is provided in <https://github.com/Div99/IQ-Learn>. For the remaining benchmarks including BC, GAIL and AIRL, we refer to a open-source implementation: https://github.com/KamyarGh/rl_swiss.

A.2. Transfer Learning Across Changing Dynamics.

In this experiment, we follow the setup in (Fu et al., 2017). A standard ant (Custom-Ant) and an ant with two disabled legs (Disabled-Ant) are simulated in Mujoco. For all benchmark algorithms tested in this experiment, we follow same network structure and hyperparameter settings described in Section A.1.

Here, we provide a supplementary experiment result to show the performance of benchmark algorithms under different number of expert trajectory. The performance of AIRL and IQ-Learn is not presented in Table 3 since we found their training is unstable (as we have mentioned in Sec.7). The scores in Hopper are recorded after 1×10^6 environment steps and the scores in other environments are recorded after 2×10^6 environment steps. The scores are reported after 3 independent Monte Carlo (MC) trials for each algorithm.

³<https://github.com/openai/spinningup>

Maximum-Likelihood Inverse Reinforcement Learning with Finite-Time Guarantees

Task	Hopper		
# Expert Trajectory	1	5	10
Expert Performance	3592.63	3530.63 ± 2.73	3531.72 ± 6.41
BC	20.53 ± 4.16	102.74 ± 63.12	371.48 ± 44.97
GAIL	2757.88 ± 293.18	2762.77 ± 209.24	3055.74 ± 179.05
f -IRL	2993.54 ± 252.59	3116.02 ± 61.30	3207.06 ± 74.67
ML-IRL(State-Only)	3138.52 ± 26.62	3131.45 ± 53.49	3164.36 ± 36.20
ML-IRL(State-Action)	3262.89 ± 24.31	3290.02 ± 62.15	3138.10 ± 144.33

Task	Half-Cheetah		
# Expert Trajectory	1	5	10
Expert Performance	5098.30	5072.53 ± 145.12	5043.02 ± 104.32
BC	-1.77 ± 0.27	155.64 ± 95.45	312.84 ± 207.99
GAIL	3254.51 ± 87.40	3085.18 ± 53.90	2963.97 ± 113.95
f -IRL	4650.49 ± 180.94	4751.63 ± 154.69	4881.13 ± 180.95
ML-IRL(State-Only)	4494.36 ± 112.03	4661.04 ± 108.05	4551.33 ± 4.72
ML-IRL(State-Action)	4728.20 ± 103.78	4846.43 ± 16.65	4894.62 ± 17.91

Task	Walker		
# Expert Trajectory	1	5	10
Expert Performance	5344.21	5471.58 ± 13.93	5471.70 ± 10.59
BC	-13.98 ± 0.22	283.425 ± 2.94	443.93 ± 63.02
GAIL	1023.41 ± 81.64	3610.49 ± 65.32	3772.54 ± 47.62
f -IRL	4483.14 ± 155.72	4562.48 ± 198.89	5028.67 ± 119.59
ML-IRL(State-Only)	4294.27 ± 55.01	4367.81 ± 162.56	4683.75 ± 9.56
ML-IRL(State-Action)	4623.46 ± 11.52	4703.35 ± 57.10	4914.42 ± 42.77

Task	Ant		
# Expert Trajectory	1	5	10
Expert Performance	5926.18	5856.84 ± 79.48	5901.73 ± 122.40
BC	760.14 ± 0.53	961.58 ± 0.17	1003.3 ± 53.93
GAIL	1107.98 ± 531.90	2971.57 ± 569.31	4008.49 ± 271.49
f -IRL	4853.53 ± 213.72	5176.09 ± 255.90	5208.41 ± 324.22
ML-IRL(State-Only)	4681.76 ± 54.35	4832.38 ± 123.28	5198.31 ± 212.81
ML-IRL(State-Action)	4919.80 ± 11.52	5157.03 ± 123.28	5275.27 ± 205.32

Task	Humanoid		
# Expert Trajectory	1	5	10
Expert Performance	5351.08	5339.12 ± 22.21	5343.76 ± 24.19
BC	76.26 ± 0.11	547.62 ± 20.21	593.64 ± 15.27
GAIL	2525.65 ± 1608.68	3174.66 ± 961.00	2966.86 ± 900.65
f -IRL	5228.58 ± 164.30	5399.67 ± 108.59	5434.54 ± 46.41
ML-IRL(State-Only)	5140.92 ± 261.28	5149.39 ± 574.27	5448.46 ± 239.75
ML-IRL(State-Action)	4328.16 ± 566.98	5281.93 ± 185.42	5290.02 ± 54.35

Table 3. Mujoco Results. The performance versus different number of expert trajectory.

B. Auxiliary Lemmas

Throughout this section, we assume Assumptions 5.1 - 5.2 hold true.

Lemma B.1. ((Xu et al., 2020, Lemma 3)) Consider the initialization distribution $\eta(\cdot)$ and transition kernel $\mathcal{P}(\cdot|s, a)$. Under $\eta(\cdot)$ and $\mathcal{P}(\cdot|s, a)$, denote $d_w(\cdot, \cdot)$ as the state-action visitation distribution of MDP with the Boltzman policy parameterized by parameter w . Suppose Assumption 5.1 holds, for all policy parameter w and w' , we have

$$\|d_w(\cdot, \cdot) - d_{w'}(\cdot, \cdot)\|_{TV} \leq C_d \|w - w'\| \quad (14)$$

where C_d is a positive constant.

Lemma B.2. ((Haarnoja et al., 2017, Theorem 4)) Under a reward function $r(\cdot, \cdot)$, given a policy π , we define a new policy $\tilde{\pi}$ as

$$\tilde{\pi}(\cdot|s) \propto \exp\left(Q_{r, \pi}^{\text{soft}}(s, \cdot)\right), \quad \forall s \in \mathcal{S}.$$

For any $s \in \mathcal{S}, a \in \mathcal{A}$, it holds that $Q_{r, \tilde{\pi}}^{\text{soft}}(s, a) \geq Q_{r, \pi}^{\text{soft}}(s, a)$.

Next, in order to facilitate analysis for entropy-regularized MDPs, we introduce a ‘‘soft’’ Bellman optimality operator $\mathcal{T} : \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ as follows:

$$\mathcal{T}(Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\max_{\pi(\cdot|s')} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a') - \log \pi(a'|s')] \right]. \quad (15)$$

In the following lemma, the properties of entropy-regularized MDPs are characterized.

Lemma B.3. ((Cen et al., 2021, Lemma 2)) The operator \mathcal{T} as defined in (15) satisfies the properties below:

- \mathcal{T} has the following closed-form expression:

$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\log \left(\sum_{a'} \exp(Q(s', a')) \right) \right]. \quad (16)$$

- \mathcal{T} is a γ -contraction in the ℓ_∞ norm, namely, for any $Q_1, Q_2 \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, it holds that

$$\|\mathcal{T}(Q_1) - \mathcal{T}(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty. \quad (17)$$

- Under a given reward function $r(\cdot, \cdot)$, the corresponding optimal soft Q -function $Q_{r, \pi^*}^{\text{soft}}$ is a unique fixed point of the operator \mathcal{T} , namely,

$$\mathcal{T}(Q_{r, \pi^*}^{\text{soft}}) = Q_{r, \pi^*}^{\text{soft}} \quad (18)$$

Proof. This Lemma is proved in (Cen et al., 2021, Lemma). We refine its analysis as below.

We first show that

$$\mathbb{E}_{a \sim \pi(\cdot|s)} \left[Q(s, a) - \log \pi(a|s) \right] = \sum_a \pi(a|s) \log \left(\frac{\exp(Q(s, a))}{\pi(a|s)} \right) \stackrel{(i)}{\leq} \log \left(\sum_a \exp(Q(s, a)) \right) \quad (19)$$

where (i) is from Jensen’s inequality. Moreover, the equality between both sides of (i) holds when the policy π has the expression $\pi(\cdot|s) \propto \exp(Q(s, \cdot))$. Therefore, through applying the inequality (19) to (15), it obtains that

$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\log \left(\sum_{a'} \exp(Q(s', a')) \right) \right], \quad (20)$$

which proves the equality (16).

We define $\|Q_1 - Q_2\|_\infty := \max_{s \in \mathcal{S}, a \in \mathcal{A}} |Q_1(s, a) - Q_2(s, a)|$ and $\epsilon = \|Q_1 - Q_2\|_\infty$. Then for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$, it follows that

$$\begin{aligned} \log \left(\sum_a \exp(Q_1(s, a)) \right) &\leq \log \left(\sum_a \exp(Q_2(s, a) + \epsilon) \right) \\ &= \log \left(\exp(\epsilon) \sum_a \exp(Q_2(s, a)) \right) \\ &= \epsilon + \log \left(\sum_a \exp(Q_2(s, a)) \right) \end{aligned}$$

Similarly, it is easy to obtain that $\log \left(\sum_a \exp(Q_1(s, a)) \right) \geq -\epsilon + \log \left(\sum_a \exp(Q_2(s, a)) \right)$. Hence, it leads to the contraction property that

$$\|\mathcal{T}(Q_1) - \mathcal{T}(Q_2)\|_\infty \leq \gamma\epsilon = \gamma\|Q_1 - Q_2\|_\infty \quad (21)$$

which proves the contraction property (17).

Moreover, we have

$$\mathcal{T}(Q_{r, \pi^*}^{\text{soft}})(s, a) \stackrel{(i)}{=} r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\log \left(\sum_{a'} \exp(Q_{r, \pi^*}^{\text{soft}}(s', a')) \right) \right] \stackrel{(ii)}{=} Q_{r, \pi^*}^{\text{soft}}(s, a) \quad (22)$$

where (i) follows the equality (20). Based on the definition of the soft Q-function $Q_{r, \pi^*}^{\text{soft}}$, we have

$$Q_{r, \pi^*}^{\text{soft}}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \left[\mathbb{E}_{a' \sim \pi^*(\cdot|s')} [-\log \pi^*(a'|s') + Q_{r, \pi^*}^{\text{soft}}(s', a')] \right]. \quad (23)$$

We prove the equality (ii) in (22) through combining (23) and the fact that the optimal soft policy has the closed form $\pi^*(\cdot|s) \propto \exp(Q_{r, \pi^*}^{\text{soft}}(s, \cdot))$. Suppose two different fixed points of the soft Bellman operator exist, then it contradicts with the contraction property in (21).

Hence, we proved the uniqueness of the optimal soft Q-function $Q_{r, \pi^*}^{\text{soft}}$. Moreover, the optimal soft Q-function $Q_{r, \pi^*}^{\text{soft}}$ is a fixed point to the soft Bellman operator \mathcal{T} in (18). \square

Lemma B.4. *Suppose Assumption 5.2 holds. Under an arbitrary policy π , for any $s \in \mathcal{S}$, $a \in \mathcal{A}$ and any reward parameters θ_1 and θ_2 , the following inequality holds:*

$$|Q_{r_{\theta_1}, \pi}^{\text{soft}}(s, a) - Q_{r_{\theta_2}, \pi}^{\text{soft}}(s, a)| \leq L_q \|\theta_1 - \theta_2\|,$$

where $L_q := \frac{L_r}{1-\gamma}$ and L_r is the positive constant in Assumption 5.2.

Proof. Based on the definition of soft-Q function, we have

$$Q_{r, \pi}^{\text{soft}}(s, a) := r(s, a) + \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \gamma^t \left(r(s_t, a_t) + \mathcal{H}(\pi(\cdot|s_t)) \right) \middle| (s_0, a_0) = (s, a) \right].$$

Then it holds that

$$\begin{aligned}
 & |Q_{r_{\theta_1}, \pi}^{\text{soft}}(s, a) - Q_{r_{\theta_2}, \pi}^{\text{soft}}(s, a)| \\
 &= \left| \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t; \theta_1) - r(s_t, a_t; \theta_2) \right) \right] \right| \\
 &\stackrel{(i)}{\leq} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left| r(s_t, a_t; \theta_1) - r(s_t, a_t; \theta_2) \right| \right] \\
 &\stackrel{(ii)}{\leq} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left\| \max_{\theta} \nabla_{\theta} r(s_t, a_t; \theta) \right\| \cdot \left\| \theta_1 - \theta_2 \right\| \right] \\
 &\stackrel{(iii)}{\leq} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t L_r \left\| \theta_1 - \theta_2 \right\| \right] \\
 &= \frac{L_r}{1 - \gamma} \left\| \theta_1 - \theta_2 \right\|
 \end{aligned} \tag{24}$$

where (i) follows Jensen's inequality; (ii) follows the mean value theorem; (iii) follows inequality (10a) in Assumption 5.2. \square

C. Proof of Lemma 4.1

Proof. First, we are able to express the objective function $L(\theta)$ in (ML-IRL) as below:

$$L(\theta) := \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \log \pi_{\theta}(a_t | s_t) \right] \stackrel{(i)}{=} \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \log \left(\frac{\exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t, a_t))}{\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t, a))} \right) \right]$$

where (i) is due to the fact that the optimal policy has the closed form $\pi_{\theta}(\cdot | s) \propto \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, \cdot))$. Therefore, we could express the objective function in this form:

$$\begin{aligned}
 L(\theta) &:= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \left(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t, a_t) - \log \left(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t, a)) \right) \right) \right] \\
 &\stackrel{(i)}{=} \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \left(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t, a_t) - V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t) \right) \right] \\
 &= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t; \theta) + \gamma V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_{t+1}) - V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t) \right) \right] \\
 &= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t; \theta) \right] + \mathbb{E}_{\tau \sim D} \left[\sum_{t=1}^{\infty} \gamma^t V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t) \right] - \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_t) \right] \\
 &= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0) \right]
 \end{aligned} \tag{25}$$

$$\stackrel{(ii)}{=} \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta} \left[\log \left(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a)) \right) \right] \tag{26}$$

where (i) and (ii) follows the fact that the the optimal soft value function could be expressed as $V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s) = \log \left(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a)) \right)$.

Based on (26), we calculate the exact gradient of the objective function $L(\theta)$ as below:

$$\begin{aligned}
 \nabla L(\theta) &:= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[\nabla_{\theta} \log \left(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a)) \right) \right] \\
 &= \mathbb{E}_{\tau \sim D} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[\sum_a \left(\frac{\exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a))}{\sum_{\bar{a}} \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, \bar{a}))} \nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a) \right) \right] \\
 &= \mathbb{E}_{\tau \sim D} \left[\sum_{k=0}^{\infty} \gamma^k \nabla_{\theta} r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[\sum_a \pi_{\theta}(a|s_0) \nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a) \right] \tag{27}
 \end{aligned}$$

Then we need to calculate the gradient $\nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a_0)$ as follows.

$$\begin{aligned}
 &\nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_0, a_0) \\
 &\stackrel{(i)}{=} \nabla_{\theta} \left(r(s_0, a_0; \theta) + \gamma \mathbb{E}_{s_1 \sim \mathcal{P}(\cdot|s_0, a_0)} \left[V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1) \right] \right) \\
 &\stackrel{(ii)}{=} \nabla_{\theta} r(s_0, a_0; \theta) + \gamma \mathbb{E}_{s_1 \sim \mathcal{P}(\cdot|s_0, a_0)} \left[\nabla_{\theta} \log \left(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1, a)) \right) \right] \\
 &= \nabla_{\theta} r(s_0, a_0; \theta) + \gamma \mathbb{E}_{s_1 \sim \mathcal{P}(\cdot|s_0, a_0)} \left[\sum_a \frac{\exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1, a))}{\sum_{\bar{a}} \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1, \bar{a}))} \nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1, a) \right] \\
 &\stackrel{(iii)}{=} \nabla_{\theta} r(s_0, a_0; \theta) + \gamma \mathbb{E}_{s_1 \sim \mathcal{P}(\cdot|s_0, a_0)} \left[\sum_a \pi_{\theta}(a|s_1) \nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_1, a) \right] \\
 &\stackrel{(iv)}{=} \nabla_{\theta} r(s_0, a_0; \theta) + \gamma \mathbb{E}_{s_1 \sim \mathcal{P}(\cdot|s_0, a_0), a_1 \sim \pi_{\theta}(\cdot|s_1)} \left[\nabla_{\theta} \left(r(s_1, a_1; \theta) + \gamma \mathbb{E}_{s_2 \sim \mathcal{P}(\cdot|s_1, a_1)} \left[V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s_2) \right] \right) \right] \\
 &\stackrel{(v)}{=} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \nabla_{\theta} r(s_t, a_t; \theta) \mid s_0, a_0 \right] \tag{28}
 \end{aligned}$$

where (i) and (iv) follows the definition of the soft Q-function; (ii) follows the fact that $V_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s) = \log(\sum_a \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a)))$; (iii) follows the fact that $\pi_{\theta}(a|s) \propto \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a))$; (v) is shown by recursively applying (i) - (iv).

Finally, plugging equation (28) into (27), the gradient of the maximum likelihood objective is:

$$\nabla L(\theta) = \mathbb{E}_{\tau \sim D} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right] - \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right]. \tag{29}$$

□

D. Proof of Lemma 5.3

To proof Lemma 5.3, we proof the equality (11a) and the equality (11b) respectively. The constants L_q and L_c in Lemma 5.3 has the expression:

$$L_q := \frac{L_r}{1 - \gamma}, \quad L_c := \frac{2L_q L_r C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}}{1 - \gamma} + \frac{2L_g}{1 - \gamma}.$$

D.1. Proof of Inequality (11a)

In this subsection, we prove the inequality (11a) in Lemma 5.3.

Proof. We show that $Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}$ has bounded gradient with respect to any reward parameter θ , then the inequality (11a) holds due to the mean value theorem. According to the equality (28), we have shown the explicit expression of $\nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a)$

for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$. Using this expression, we have the following series of relations:

$$\begin{aligned}
 \|\nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a)\| &\stackrel{(i)}{=} \left\| \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \mid (s_0, a_0) = (s, a) \right] \right\| \\
 &\stackrel{(ii)}{\leq} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t \left\| \nabla_{\theta} r(s_t, a_t; \theta) \right\| \mid (s_0, a_0) = (s, a) \right] \\
 &\stackrel{(iii)}{\leq} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t L_r \mid (s_0, a_0) = (s, a) \right] \\
 &= \frac{L_r}{1 - \gamma}
 \end{aligned} \tag{30}$$

where (i) is from the equality (28) in the proof of Lemma 4.1, (ii) follows Jensen's inequality and (iii) follows the inequality (10a) in Assumption 5.2. To complete this proof, we use the mean value theorem to show that

$$|Q_{r_{\theta_1}, \pi_{\theta_1}}^{\text{soft}}(s, a) - Q_{r_{\theta_2}, \pi_{\theta_2}}^{\text{soft}}(s, a)| \leq \|\max_{\theta} \nabla_{\theta} Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a)\| \cdot \|\theta_1 - \theta_2\| \leq L_q \|\theta_1 - \theta_2\|$$

where the last inequality follows (30) and we denote $L_q := \frac{L_r}{1 - \gamma}$. Therefore, we have proved the Lipschitz continuous inequality in (11a). \square

D.2. Proof of Inequality (11b)

In this section, we prove the inequality (11b) in Lemma 5.3.

Proof. According to Lemma 4.1, the gradient $\nabla L(\theta)$ is expressed as:

$$\nabla L(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right] - \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) \right]. \tag{31}$$

Using the above relation, we have

$$\begin{aligned}
 &\|\nabla L(\theta_1) - \nabla L(\theta_2)\| \\
 &\stackrel{(i)}{=} \left\| \left(\mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] \right) - \right. \\
 &\quad \left. \left(\mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right) \right\| \\
 &\leq \underbrace{\left\| \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right\|}_{:= \text{term A}} + \\
 &\quad \underbrace{\left\| \mathbb{E}_{\tau \sim \pi_{\theta_1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right\|}_{:= \text{term B}}
 \end{aligned} \tag{32}$$

where (i) follows the exact gradient expression in equation (31). Then we separately analyze term A and term B in (32).

For term A, it follows that

$$\begin{aligned}
 & \left\| \mathbb{E}_{\tau \sim D} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim D} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right\| \\
 & \stackrel{(i)}{\leq} \mathbb{E}_{\tau \sim D} \left[\sum_{t \geq 0} \gamma^t \left\| \nabla_{\theta} r(s_t, a_t; \theta_1) - \nabla_{\theta} r(s_t, a_t; \theta_2) \right\| \right] \\
 & \stackrel{(ii)}{\leq} \mathbb{E}_{\tau \sim D} \left[\sum_{t \geq 0} \gamma^t L_g \|\theta_1 - \theta_2\| \right] \\
 & = \frac{L_g}{1 - \gamma} \|\theta_1 - \theta_2\|
 \end{aligned} \tag{33}$$

where (i) follows Jensen's inequality and (ii) is from (10b) in Assumption 5.2.

For the term B, it holds that

$$\begin{aligned}
 & \left\| \mathbb{E}_{\tau \sim \pi_{\theta_1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right\| \\
 & \stackrel{(i)}{\leq} \left\| \mathbb{E}_{\tau \sim \pi_{\theta_1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] \right\| \\
 & \quad + \left\| \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_2) \right] \right\| \\
 & \stackrel{(ii)}{\leq} \frac{1}{1 - \gamma} \left\| \mathbb{E}_{(s,a) \sim d(\cdot, \cdot; \pi_{\theta_1})} \left[\nabla_{\theta} r(s_t, a_t; \theta_1) \right] - \mathbb{E}_{(s,a) \sim d(\cdot, \cdot; \pi_{\theta_2})} \left[\nabla_{\theta} r(s_t, a_t; \theta_1) \right] \right\| \\
 & \quad + \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{t \geq 0} \gamma^t \left\| \nabla_{\theta} r(s_t, a_t; \theta_1) - \nabla_{\theta} r(s_t, a_t; \theta_2) \right\| \right] \\
 & \stackrel{(iii)}{\leq} \frac{1}{1 - \gamma} \left\| \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \nabla_{\theta} r(s_t, a_t; \theta_1) \left(d(s, a; \pi_{\theta_1}) - d(s, a; \pi_{\theta_2}) \right) \right\| + \mathbb{E}_{\tau \sim \pi_{\theta_2}} \left[\sum_{k \geq 0} \gamma^k L_g \|\theta_1 - \theta_2\| \right] \\
 & \stackrel{(iv)}{\leq} \frac{2L_r}{1 - \gamma} \|d(\cdot, \cdot; \pi_{\theta_1}) - d(\cdot, \cdot; \pi_{\theta_2})\|_{TV} + \frac{L_g}{1 - \gamma} \|\theta_1 - \theta_2\|
 \end{aligned} \tag{34}$$

where (i) follows the triangle inequality, (ii) is from Jensen's inequality and the definition of the discounted state-action visitation measure $d(s, a; \pi) := (1 - \gamma) \pi(a|s) \sum_{t \geq 0} \gamma^t \mathcal{P}^{\pi}(s_t = s)$; (iii) is from (10b) in Assumption 5.2; (iv) is from (10a) and the definition of the total variation norm.

Plugging the inequalities (33), (34) to (32), it holds that

$$\begin{aligned}
 & \|\nabla L(\theta_1) - \nabla L(\theta_2)\| \\
 & \leq \frac{2L_r}{1 - \gamma} \|d(\cdot, \cdot; \pi_{\theta_1}) - d(\cdot, \cdot; \pi_{\theta_2})\|_{TV} + \frac{2L_g}{1 - \gamma} \|\theta_1 - \theta_2\| \\
 & \stackrel{(i)}{\leq} \frac{2L_r C_d}{1 - \gamma} \|Q_{r_{\theta_1}, \pi_{\theta_1}}^{\text{soft}} - Q_{r_{\theta_2}, \pi_{\theta_2}}^{\text{soft}}\| + \frac{2L_g}{1 - \gamma} \|\theta_1 - \theta_2\| \\
 & \stackrel{(ii)}{\leq} \frac{2L_r C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}}{1 - \gamma} \|Q_{r_{\theta_1}, \pi_{\theta_1}}^{\text{soft}} - Q_{r_{\theta_2}, \pi_{\theta_2}}^{\text{soft}}\|_{\infty} + \frac{2L_g}{1 - \gamma} \|\theta_1 - \theta_2\| \\
 & \stackrel{(iii)}{\leq} \left(\frac{2L_q L_r C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}}{1 - \gamma} + \frac{2L_g}{1 - \gamma} \right) \|\theta_1 - \theta_2\|.
 \end{aligned} \tag{35}$$

Given the fact that π_{θ} is a Boltzmann policy parameterized by $Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}$ where $\pi_{\theta}(a|s) \propto \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a))$, we show the inequality (i) from the inequality (14) in Lemma B.1. Moreover, the inequality (ii) follows the equivalence relation between Frobenius norm and infinity norm and (iii) is from the inequality (11a) in Lemma 5.3.

Define the constant $L_c := \frac{2L_q L_r C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}}{1-\gamma} + \frac{2L_g}{1-\gamma}$, we have the following inequality:

$$\|\nabla L(\theta_1) - \nabla L(\theta_2)\| \leq L_c \|\theta_1 - \theta_2\|.$$

Therefore, we complete the proof of the inequality (11b) in Lemma 5.3. \square

E. Proof of Lemma 6.1

Proof. Suppose the expert trajectories τ in **ML-IRL** is sampled from an expert policy π^E . Moreover, we parameterize the state-only reward as $r(s; \theta)$. Then the objective function $L(\theta)$ in **ML-IRL** could be rewritten as follows.

$$\begin{aligned} L(\theta) &:= \mathbb{E}_{\tau \sim \pi^E} \left[\sum_{t \geq 0} \gamma^t \log \pi_\theta(a_t | s_t) \right] \\ &\stackrel{(i)}{=} \mathbb{E}_{\tau \sim \pi^E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] \\ &\stackrel{(ii)}{=} \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi^E}^{\text{soft}}(s_0) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] - H(\pi^E) \end{aligned} \quad (36)$$

where (i) follows (25) and the fact that the reward is a state-only function $r(s; \theta)$; (ii) follows the definitions of the soft value function.

Ignoring the constant term $H(\pi^E)$ in (36), the maximum likelihood formulation (**ML-IRL**) is equivalent to the following bi-level problem:

$$\begin{aligned} \min_{\theta} \quad & \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi^E}^{\text{soft}}(s_0) \right] \\ \text{s.t. } \quad & \pi_\theta := \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t; \theta) + \mathcal{H}(\pi(\cdot | s_t)) \right) \right]. \end{aligned}$$

Therefore, we complete the proof of Lemma 6.1. As an alternative interpretation to (**ML-IRL**), the formulation above aims to minimize the gap between the soft value function of π_θ and π^E under the state-only IRL setting. \square

F. Proof of Theorem 3.1

Proof. Calculate the Lagrangian of MaxEnt-IRL, we obtain that

$$\begin{aligned}
 H(\pi) &+ \left\langle \theta, \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle + \sum_{s \in \mathcal{S}, t \geq 0} C_{s_t=s} \left(1 - \sum_{a \in \mathcal{A}} \pi(a|s_t) \right) \\
 &= \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t = s) \right] + \left\langle \theta, \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle \\
 &+ \sum_{s \in \mathcal{S}, t \geq 0} C_{s_t=s} \left(1 - \sum_{a \in \mathcal{A}} \pi(a|s_t = s) \right)
 \end{aligned} \tag{37}$$

where θ is the dual variable to ensure the feature matching equality in (1), and $C_{s_t=s}$ is the dual variable to ensure that π is a well-defined policy satisfying $\sum_{a \in \mathcal{A}} \pi(a|s_t = s) = 1$.

Then we could calculate the gradient of (37) w.r.t. $\pi(a|s_t = s)$, and set it to 0. Then it holds that

$$\begin{aligned}
 0 &= \mathcal{P}^\pi(s_t = s) \left(-\gamma^t (\log \pi(a|s_t = s) + 1) + \mathbb{E}_\pi \left[\sum_{\kappa=t}^{\infty} -\gamma^{\kappa+1} \log \pi(a_{\kappa+1} | s_{\kappa+1}) \mid s_t = s, a_t = a \right] \right. \\
 &\quad \left. + \theta^T \mathbb{E}_\pi \left[\sum_{\kappa=t}^{\infty} \gamma^\kappa \phi(s_\kappa, a_\kappa) \mid s_t = s, a_t = a \right] \right) - C_{s_t=s}.
 \end{aligned} \tag{38}$$

Dividing $\gamma^t \mathcal{P}^\pi(s_t = s)$ on both sides of (38) and further moving $\log \pi(a|s_t = s)$ to the left side, then we have the equality as below:

$$\begin{aligned}
 \log \pi(a|s_t = s) &= \left(-\frac{C_{s_t=s}}{\gamma^t \mathcal{P}^\pi(s_t = s)} - 1 \right) + \mathbb{E}_\pi \left[\sum_{\kappa=t+1}^{\infty} -\gamma^{\kappa-t} \log \pi(a_\kappa | s_\kappa) \mid s_t = s, a_t = a \right] \\
 &\quad + \theta^T \mathbb{E}_\pi \left[\sum_{\kappa=t}^{\infty} \gamma^{\kappa-t} \phi(s_\kappa, a_\kappa) \mid s_t = s, a_t = a \right]
 \end{aligned} \tag{39}$$

Given that $-\frac{C_{s_t=s}}{\gamma^t \mathcal{P}^\pi(s_t = s)} - 1$ is independent of action a , we could express the closed form of $\pi(a|s_t = s)$ as below:

$$\pi(a|s_t = s) \propto \exp \left(\mathbb{E}_\pi \left[\sum_{\kappa=t+1}^{\infty} -\gamma^{\kappa-t} \log \pi(a_\kappa | s_\kappa) \mid s_t = s, a_t = a \right] + \theta^T \mathbb{E}_\pi \left[\sum_{\kappa=t}^{\infty} \gamma^{\kappa-t} \phi(s_\kappa, a_\kappa) \mid s_t = s, a_t = a \right] \right).$$

According to the closed form of the policy above, it shows that $\pi(a|s_t = s)$ is a stationary policy being independent of the time index t . Therefore, it holds that $\pi(a|s_t = s) = \pi(a|s)$ for any $t \geq 0$.

Denoting a linearly parameterized reward as $r(s, a; \theta) := \theta^T \phi(s, a)$, it holds that

$$\begin{aligned}
 \pi(a|s) &\propto \exp \left(\theta^T \mathbb{E}_\pi \left[\sum_{\kappa=0}^{\infty} \gamma^\kappa \phi(s_\kappa, a_\kappa) \mid s_0 = s, a_0 = a \right] + \mathbb{E}_\pi \left[\sum_{\kappa=0}^{\infty} -\gamma^{\kappa+1} \log \pi(a_{\kappa+1} | s_{\kappa+1}) \mid s_0 = s, a_0 = a \right] \right) \\
 &= \exp \left(\mathbb{E}_\pi \left[\sum_{\kappa=0}^{\infty} \gamma^\kappa r(s_\kappa, a_\kappa; \theta) \mid s_0 = s, a_0 = a \right] + \mathbb{E}_\pi \left[\sum_{\kappa=0}^{\infty} -\gamma^{\kappa+1} \log \pi(a_{\kappa+1} | s_{\kappa+1}) \mid s_0 = s, a_0 = a \right] \right)
 \end{aligned} \tag{40}$$

Here, the optimal $\pi(a|s)$ is a function of the dual variables (reward parameters) θ . In the maximum entropy reinforcement learning (Haarnoja et al., 2017), under a reward function $r(\cdot, \cdot)$ and policy π , the soft value function and soft Q-function are defined as below:

$$V_{r, \pi}^{\text{soft}}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t) + \mathcal{H}(\pi(\cdot | s_t)) \right) \mid s_0 = s \right] \tag{41a}$$

$$Q_{r, \pi}^{\text{soft}}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [V_{r, \pi}^{\text{soft}}(s')] \tag{41b}$$

Based on the definitions in (41a) - (41b), we could further express the closed form of the policy in (40) as below:

$$\pi(a|s) = \frac{\exp(Q_{r_\theta, \pi}^{\text{soft}}(s, a))}{\sum_{a \in \mathcal{A}} \exp(Q_{r_\theta, \pi}^{\text{soft}}(s, a))}. \quad (42)$$

According to (Haarnoja et al., 2017), under a reward function $r(\cdot, \cdot)$, the optimal soft policy π satisfies $\pi(\cdot|s) \propto \exp(Q_{r, \pi}^{\text{soft}}(s, \cdot))$. Hence, we have shown that the policy in (42) is the optimal policy under the reward function $r(\cdot, \cdot; \theta)$. After denoting the optimal policy under $r(\cdot, \cdot; \theta)$ as π_θ , we have the following relation:

$$\pi_\theta(a|s) = \frac{\exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))}{\sum_{a \in \mathcal{A}} \exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))} \stackrel{(a)}{=} \frac{\exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))}{\exp(V_{r_\theta, \pi_\theta}^{\text{soft}}(s))} = \exp\left(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a) - V_{r_\theta, \pi_\theta}^{\text{soft}}(s)\right) \quad (43)$$

where (a) is due to the equality shown as below:

$$\begin{aligned} V_{r_\theta, \pi_\theta}^{\text{soft}}(s) &= \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[-\log(\pi_\theta(a|s)) + Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a) \right] \\ &= \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[-\log\left(\frac{\exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))}{\sum_{a \in \mathcal{A}} \exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))}\right) + Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a) \right] \\ &= \log\left(\sum_{a \in \mathcal{A}} \exp(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a))\right). \end{aligned}$$

Rewriting the equality (38), we are able to show the expression of $C_{s_t=s}$ as below:

$$\begin{aligned} C_{s_t=s} &= \mathcal{P}^\pi(s_t = s) \left(-\gamma^t (\log \pi_\theta(a|s_t = s) + 1) + \mathbb{E}_{\pi_\theta} \left[\sum_{\kappa=t}^{\infty} -\gamma^{\kappa+1} \log \pi_\theta(a_{\kappa+1}|s_{\kappa+1}) \mid s_t = s, a_t = a \right] \right. \\ &\quad \left. + \theta^T \mathbb{E}_{\pi_\theta} \left[\sum_{\kappa=t}^{\infty} \gamma^\kappa \phi(s_\kappa, a_\kappa) \mid s_t = s, a_t = a \right] \right) \\ &= \gamma^t \mathcal{P}^{\pi_\theta}(s_t = s) \left(-1 - \log \pi_\theta(a|s) + \mathbb{E}_{\pi_\theta} \left[\sum_{\kappa=0}^{\infty} -\gamma^{\kappa+1} \log \pi_\theta(a_{\kappa+1}|s_{\kappa+1}) \mid s_0 = s, a_0 = a \right] \right. \\ &\quad \left. + \mathbb{E}_{\pi_\theta} \left[\sum_{\kappa=0}^{\infty} \gamma^\kappa r(s_\kappa, a_\kappa; \theta) \mid s_0 = s, a_0 = a \right] \right) \\ &\stackrel{(a)}{=} \gamma^t \mathcal{P}^\pi(s_t = s) \left(-1 - \log \pi_\theta(a|s) + Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a) \right) \\ &\stackrel{(b)}{=} \gamma^t \mathcal{P}^\pi(s_t = s) \left(V_{r_\theta, \pi_\theta}^{\text{soft}}(s) - 1 \right) \end{aligned} \quad (44)$$

where (a) follows the definition of the soft Q-function in (41b), and (b) follows (43). According to (44), we are able to show the exact expression of $C_{s_t=s}$.

Plugging π_θ and $C_{s_t=s}$ into (37), we have

$$\begin{aligned}
 & \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi_\theta(a_t | s_t) \right] + \left\langle \theta, \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle \\
 & + \sum_{s \in \mathcal{S}, t \geq 0} C_{s_t=s} \left(1 - \sum_{a \in \mathcal{A}} \pi_\theta(a | s_t = s) \right) \\
 & \stackrel{(a)}{=} \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi_\theta(a_t | s_t) \right] + \left\langle \theta, \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle \\
 & \stackrel{(b)}{=} \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} -\gamma^t \left(Q_{r_\theta, \pi_\theta}^{\text{soft}}(s, a) - V_{r_\theta, \pi_\theta}^{\text{soft}}(s) \right) \right] + \left\langle \theta, \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle \\
 & \stackrel{(c)}{=} \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} -\gamma^t \left(\theta^T \phi(s_t, a_t) + \gamma V_{r_\theta, \pi_\theta}^{\text{soft}}(s_{t+1}) - V_{r_\theta, \pi_\theta}^{\text{soft}}(s_t) \right) \right] \\
 & + \left\langle \theta, \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \right\rangle \\
 & = \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] - \theta^T \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] \tag{45}
 \end{aligned}$$

where (a) is due to the fact that $\sum_{a \in \mathcal{A}} \pi_\theta(a | s_t = s) = 1$ for all $a \in \mathcal{A}$ and $s \in \mathcal{S}$; (b) follows (43); (c) is due to the definition of the soft Q-function in (41b). Here, we could further show the problem in (45) is equivalent to ML-IRL as below:

$$\begin{aligned}
 \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \ln \pi_\theta(a_t | s_t) \right] &= \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{\tau \sim \mathcal{D}} \left[r(s_t, a_t; \theta) + \gamma V_{r_\theta, \pi_\theta}^{\text{soft}}(s_{t+1}) - V_{r_\theta, \pi_\theta}^{\text{soft}}(s_t) \right] \\
 &= \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t; \theta) \right] + \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{\tau \sim \mathcal{D}} \left[\gamma V_{r_\theta, \pi_\theta}^{\text{soft}}(s_{t+1}) - V_{r_\theta, \pi_\theta}^{\text{soft}}(s_t) \right] \\
 &= \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t; \theta) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] \\
 &= \theta^T \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t, a_t) \right] - \mathbb{E}_{s_0 \sim \eta(\cdot)} \left[V_{r_\theta, \pi_\theta}^{\text{soft}}(s_0) \right] \tag{46}
 \end{aligned}$$

Finally, through combining (45) and (46), we are able to know that the maximum likelihood formulation ML-IRL is the dual form of MaxEnt-IRL.

□

G. Proof of Theorem 5.4

In this section, we prove (12a) and (12b) respectively, to show the convergence of the lower-level problem and the upper-level problem.

G.1. Proof of (12a)

Proof. In this proof, we first show the convergence of the lower-level variable $\{\pi_k\}_{k \geq 0}$. Recall that we approximate the optimal policy π_{θ_k} by π_{k+1} at each iteration k . We first analyze the approximation error between π_{θ_k} and π_{k+1} as follows. For any $s \in \mathcal{S}$ and $a \in \mathcal{A}$, we have the following relation:

$$\begin{aligned} & \left| \log(\pi_{k+1}(a|s)) - \log(\pi_{\theta_k}(a|s)) \right| \\ & \stackrel{(i)}{=} \left| \log\left(\frac{\exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a))}{\sum_{\tilde{a}} \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \tilde{a}))}\right) - \log\left(\frac{\exp(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a))}{\sum_{\tilde{a}} \exp(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, \tilde{a}))}\right) \right| \\ & \stackrel{(ii)}{\leq} \left| Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a) - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a) \right| + \left| \log\left(\sum_{\tilde{a}} \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \tilde{a}))\right) - \log\left(\sum_{\tilde{a}} \exp(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, \tilde{a}))\right) \right| \end{aligned} \quad (47)$$

where (i) follows (6) and the fact that $\pi_{\theta}(a|s) \propto \exp(Q_{r_{\theta}, \pi_{\theta}}^{\text{soft}}(s, a))$; (ii) follows the triangle inequality. We further analyze the second term in (47).

We first denote the operator $\log(\|\exp(v)\|_1) := \log(\|\sum_{\tilde{a} \in \mathcal{A}} \exp(v_{\tilde{a}})\|_1)$, where the vector $v \in \mathbb{R}^{|\mathcal{A}|}$ and $v = [v_1, v_2, \dots, v_{|\mathcal{A}|}]$. Then for any $v', v'' \in \mathbb{R}^{|\mathcal{A}|}$, we have the following relation:

$$\begin{aligned} \left| \log(\|\exp(v')\|_1) - \log(\|\exp(v'')\|_1) \right| & \stackrel{(i)}{=} \langle v' - v'', \nabla_v \log(\|\exp(v)\|_1)|_{v=v^c} \rangle \\ & \leq \|v' - v''\|_{\infty} \cdot \|\nabla_v \log(\|\exp(v)\|_1)|_{v=v^c}\|_1 \\ & \stackrel{(ii)}{=} \|v' - v''\|_{\infty} \end{aligned} \quad (48)$$

where (i) follows the mean value theorem and v_c is a convex combination of v' and v'' ; (ii) follows the following equalities:

$$[\nabla_v \log(\|\exp(v)\|_1)]_i = \frac{\exp(v_i)}{\sum_{1 \leq a \leq |\mathcal{A}|} \exp(v_a)}, \quad \|\nabla_v \log(\|\exp(v)\|_1)\|_1 = 1, \quad \forall v \in \mathbb{R}^{|\mathcal{A}|}.$$

Through plugging (48) into (47), it holds that

$$\begin{aligned} & \left| \log(\pi_{k+1}(a|s)) - \log(\pi_{\theta_k}(a|s)) \right| \\ & \leq \left| Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a) - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a) \right| + \max_{\tilde{a} \in \mathcal{A}} \left| Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \tilde{a}) - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, \tilde{a}) \right| \end{aligned} \quad (49)$$

Taking the infinity norm over $\mathbb{R}^{|\mathcal{S}| \cdot |\mathcal{A}|}$, the following result holds:

$$\|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty} \leq 2 \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \quad (50)$$

where $\|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty} = \max_{s \in \mathcal{S}, a \in \mathcal{A}} |\log \pi_{k+1}(a|s) - \log \pi_{\theta_k}(a|s)|$ and $\|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} = \max_{s \in \mathcal{S}, a \in \mathcal{A}} |Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a) - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a)|$.

Based on the inequality (50), we analyze $\|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty}$ to show the convergence of the policy estimates. It leads to the following analysis:

$$\begin{aligned} & \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \\ & = \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} + Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}} + Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}}\|_{\infty} \\ & \leq \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + \|Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}}\|_{\infty} \\ & \stackrel{(i)}{\leq} L_q \|\theta_k - \theta_{k-1}\| + \|Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}}\|_{\infty} \\ & \stackrel{(ii)}{\leq} \|Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + 2L_q \|\theta_k - \theta_{k-1}\| \end{aligned} \quad (51)$$

where (i) is from (11a) in Lemma 5.3; (ii) follows Lemma B.4. Based on (51), we further analyze the two terms in (51) as below.

Recall Lemma B.3, we have the ‘‘soft’’ Bellman operator expressed as below:

$$\mathcal{T}_\theta(Q)(s, a) = r(s, a; \theta) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \left[\log \left(\sum_{a'} \exp(Q(s', a')) \right) \right] \quad (52)$$

According to the soft Bellman operator, it holds that

$$\begin{aligned} Q_{r_{\theta_k}, \pi_{k+1}}^{\text{soft}}(s, a) &= r(s, a; \theta_k) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [V_{r_{\theta_k}, \pi_{k+1}}^{\text{soft}}(s')] \\ &= r(s, a; \theta_k) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a), a' \sim \pi_{k+1}(\cdot | s')} [-\log \pi_{k+1}(a' | s') + Q_{r_{\theta_k}, \pi_{k+1}}^{\text{soft}}(s', a')] \\ &\stackrel{(i)}{\geq} r(s, a; \theta_k) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a), a' \sim \pi_{k+1}(\cdot | s')} [-\log \pi_{k+1}(a' | s') + Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s', a')] \\ &\stackrel{(ii)}{=} r(s, a; \theta_k) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \left[\log \left(\sum_{a'} \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s', a')) \right) \right] \\ &\stackrel{(iii)}{=} \mathcal{T}_{\theta_k}(Q_{r_{\theta_k}, \pi_k}^{\text{soft}})(s, a) \end{aligned} \quad (53)$$

where (i) follows the policy improvement result in Lemma B.2, (ii) follows the definition $\pi_{k+1}(a | s) := \frac{\exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, a))}{\sum_{\hat{a}} \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \hat{a}))}$ in (6); (iii) follows the definition of the soft Bellman operator in (52).

For any $s \in \mathcal{S}$ and $a \in \mathcal{A}$, it holds that

$$0 \stackrel{(i)}{\leq} Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a) - Q_{r_{\theta_k}, \pi_{k+1}}^{\text{soft}}(s, a) \stackrel{(ii)}{\leq} Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, a) - \mathcal{T}_{\theta_k}(Q_{r_{\theta_k}, \pi_k}^{\text{soft}})(s, a) \quad (54)$$

where (i) is due to the fact that π_{θ_k} is the optimal policy under reward parameter θ_k ; (ii) is from (53).

Hence, it further leads to

$$\begin{aligned} \|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{k+1}}^{\text{soft}}\|_\infty &\stackrel{(i)}{\leq} \|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - \mathcal{T}_{\theta_k}(Q_{r_{\theta_k}, \pi_k}^{\text{soft}})\|_\infty \\ &\stackrel{(ii)}{=} \|\mathcal{T}_{\theta_k}(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}) - \mathcal{T}_{\theta_k}(Q_{r_{\theta_k}, \pi_k}^{\text{soft}})\|_\infty \\ &\stackrel{(iii)}{\leq} \gamma \|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|_\infty \end{aligned} \quad (55)$$

where (i) is from (54); (ii) is from the fixed-point property in (18); (iii) is from the contraction property in (17). Therefore, we have the following result:

$$\begin{aligned} &\|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_\infty \\ &\stackrel{(i)}{\leq} \|Q_{r_{\theta_{k-1}}, \pi_k}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_\infty + 2L_q \|\theta_k - \theta_{k-1}\| \\ &\stackrel{(ii)}{\leq} \gamma \|Q_{r_{\theta_{k-1}}, \pi_{k-1}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_\infty + 2L_q \|\theta_k - \theta_{k-1}\| \end{aligned} \quad (56)$$

where (i) is from (51); (ii) is from (55).

To show the convergence of the soft Q-function based on (56), we further analyze the error between the reward parameters θ_k and θ_{k-1} . Recall in Alg.1, the updates in reward parameters follows (9):

$$\begin{aligned} \theta_k &= \theta_{k-1} + \alpha g_{k-1} \\ &= \theta_{k-1} + \alpha (h(\theta_{k-1}, \tau_{k-1}^E) - h(\theta_{k-1}, \tau_{k-1}^A)) \end{aligned}$$

where we denote $\tau = \{(s_t, a_t)\}_{t=0}^{\infty}$, $h(\theta, \tau) := \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta)$ and g_{k-1} is the stochastic gradient estimator at iteration $k-1$. Here, τ_{k-1}^E denotes the trajectory sampled from the expert's dataset D at iteration $k-1$ and τ_{k-1}^A denotes the trajectory sampled from the agent's policy π_k at time $k-1$. Then according to the inequality (10a) in Assumption 5.2, we could show that

$$\|g_{k-1}\| \leq \|h(\theta_{k-1}, \tau_{k-1}^E)\| + \|h(\theta_{k-1}, \tau_{k-1}^A)\| \leq 2L_r \sum_{t \geq 0} \gamma^t = \frac{2L_r}{1-\gamma} = 2L_q \quad (57)$$

where the last equality follows the fact that we have defined the constant $L_q := \frac{L_r}{1-\gamma}$. Then we could further show that

$$\begin{aligned} & \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \\ & \stackrel{(i)}{\leq} \gamma \|Q_{r_{\theta_{k-1}}, \pi_{k-1}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + 2L_q \|\theta_k - \theta_{k-1}\| \\ & \stackrel{(ii)}{=} \gamma \|Q_{r_{\theta_{k-1}}, \pi_{k-1}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + 2\alpha L_q \|g_{k-1}\| \\ & \stackrel{(iii)}{\leq} \gamma \|Q_{r_{\theta_{k-1}}, \pi_{k-1}}^{\text{soft}} - Q_{r_{\theta_{k-1}}, \pi_{\theta_{k-1}}}^{\text{soft}}\|_{\infty} + 4\alpha L_q^2 \end{aligned} \quad (58)$$

where (i) is from (56); (ii) follows the reward update scheme in (9); (iii) is from (57).

Summing the inequality (58) from $k=1$ to $k=K$, it holds that

$$\sum_{k=1}^K \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \leq \gamma \sum_{k=0}^{K-1} \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} + 4\alpha K L_q^2 \quad (59)$$

Rearranging the inequality (59) and divided (59) by K on both sides, it holds that

$$\frac{1-\gamma}{K} \sum_{k=1}^K \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \leq \frac{\gamma}{K} \left(\|Q_{r_{\theta_0}, \pi_0}^{\text{soft}} - Q_{r_{\theta_0}, \pi_{\theta_0}}^{\text{soft}}\|_{\infty} - \|Q_{r_{\theta_K}, \pi_K}^{\text{soft}} - Q_{r_{\theta_K}, \pi_{\theta_K}}^{\text{soft}}\|_{\infty} \right) + 4\alpha L_q^2 \quad (60)$$

Dividing the constant $1-\gamma$ on both sides of (60), it holds that

$$\frac{1}{K} \sum_{k=1}^K \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \leq \frac{\gamma C_0}{K(1-\gamma)} + \frac{4L_q^2}{1-\gamma} \alpha$$

where we denote $C_0 := \|Q_{r_{\theta_0}, \pi_0}^{\text{soft}} - Q_{r_{\theta_0}, \pi_{\theta_0}}^{\text{soft}}\|_{\infty}$. We could also write the inequality above as

$$\begin{aligned} & \frac{1}{K} \sum_{k=0}^{K-1} \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} \\ & \leq \frac{\gamma C_0}{T(1-\gamma)} + \frac{C_0}{T} - \frac{\|Q_{r_{\theta_K}, \pi_K}^{\text{soft}} - Q_{r_{\theta_K}, \pi_{\theta_K}}^{\text{soft}}\|_{\infty}}{K} + \frac{4L_q^2}{1-\gamma} \alpha \\ & \leq \frac{C_0}{T(1-\gamma)} + \frac{4L_q^2}{1-\gamma} \alpha. \end{aligned}$$

Recall the stepsize is defined as $\alpha = \frac{\alpha_0}{T^\sigma}$ where $\sigma > 0$. Then we have the following result:

$$\frac{1}{K} \sum_{k=0}^{K-1} \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} = \mathcal{O}(K^{-1}) + \mathcal{O}(K^{-\sigma}). \quad (61)$$

With the inequality (50), it follows that

$$\frac{1}{K} \sum_{k=0}^{K-1} \|\log \pi_{k+1} - \log \pi_{\theta_k}\|_{\infty} \leq \frac{2}{K} \sum_{k=0}^{K-1} \|Q_{r_{\theta_k}, \pi_k}^{\text{soft}} - Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}\|_{\infty} = \mathcal{O}(K^{-1}) + \mathcal{O}(K^{-\sigma}).$$

Therefore, we complete the proof of (12a) in Theorem 5.4. \square

G.2. Proof of (12b)

Proof. In this part, we prove the convergence of reward parameters $\{\theta_k\}_{k \geq 0}$.

We have the following result of the objective function $L(\theta)$:

$$\begin{aligned}
 L(\theta_{k+1}) &\stackrel{(i)}{\geq} L(\theta_k) + \langle \nabla L(\theta_k), \theta_{k+1} - \theta_k \rangle - \frac{L_c}{2} \|\theta_{k+1} - \theta_k\|^2 \\
 &\stackrel{(ii)}{=} L(\theta_k) + \alpha \langle \nabla L(\theta_k), g_k \rangle - \frac{L_c \alpha^2}{2} \|g_k\|^2 \\
 &= L(\theta_k) + \alpha \langle \nabla L(\theta_k), g_k - \nabla L(\theta_k) \rangle + \alpha \|\nabla L(\theta_k)\|^2 - \frac{L_c \alpha^2}{2} \|g_k\|^2 \\
 &\stackrel{(iii)}{\geq} L(\theta_k) + \alpha \langle \nabla L(\theta_k), g_k - \nabla L(\theta_k) \rangle + \alpha \|\nabla L(\theta_k)\|^2 - 2L_c L_q^2 \alpha^2
 \end{aligned} \tag{62}$$

where (i) is from the Lipschitz smooth property in (11b) of Lemma 5.3; (ii) follows the update scheme (9); (iii) is from constant bound in (57).

Taking an expectation over the both sides of (62), it holds that

$$\begin{aligned}
 &\mathbb{E}[L(\theta_{k+1})] \\
 &\geq \mathbb{E}[L(\theta_k)] + \alpha \mathbb{E}[\langle \nabla L(\theta_k), g_k - \nabla L(\theta_k) \rangle] + \alpha \mathbb{E}[\|\nabla L(\theta_k)\|^2] - 2L_c L_q^2 \alpha^2 \\
 &= \mathbb{E}[L(\theta_k)] + \alpha \mathbb{E}[\langle \nabla L(\theta_k), \mathbb{E}[g_k - \nabla L(\theta_k) | \theta_k] \rangle] + \alpha \mathbb{E}[\|\nabla L(\theta_k)\|^2] - 2L_c L_q^2 \alpha^2 \\
 &\stackrel{(i)}{=} \mathbb{E}[L(\theta_k)] + \alpha \mathbb{E}[\langle \nabla L(\theta_k), \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_t) \right] - \mathbb{E}_{\tau \sim \pi_{k+1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_t) \right] \rangle] \\
 &\quad + \alpha \mathbb{E}[\|\nabla L(\theta_k)\|^2] - 2L_c L_q^2 \alpha^2 \\
 &\stackrel{(ii)}{\geq} \mathbb{E}[L(\theta_k)] - 2\alpha L_q \underbrace{\mathbb{E} \left[\left\| \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] - \mathbb{E}_{\tau \sim \pi_{k+1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] \right\| \right]}_{\text{term A}} \\
 &\quad + \alpha \mathbb{E}[\|\nabla L(\theta_k)\|^2] - 2L_c L_q^2 \alpha^2
 \end{aligned} \tag{63}$$

where (i) follows (7) and (8); (ii) is due to the fact that $\|\nabla L(\theta)\| \leq 2L_q$.

Then we further analyze the term A as below:

$$\begin{aligned}
 &\mathbb{E} \left[\left\| \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] - \mathbb{E}_{\tau \sim \pi_{k+1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] \right\| \right] \\
 &\stackrel{(i)}{=} \mathbb{E} \left[\left\| \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d(\cdot, \cdot; \pi_{\theta_k})} [\nabla_{\theta} r(s, a; \theta_k)] - \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d(\cdot, \cdot; \pi_{k+1})} [\nabla_{\theta} r(s, a; \theta_k)] \right\| \right] \\
 &\stackrel{(ii)}{\leq} \frac{2}{1-\gamma} \cdot \max_{s \in \mathcal{S}, a \in \mathcal{A}} \|\nabla_{\theta} r(s, a; \theta_k)\| \cdot \mathbb{E}[\|d(\cdot, \cdot; \pi_{\theta_k}) - d(\cdot, \cdot; \pi_{k+1})\|_{TV}] \\
 &\stackrel{(iii)}{\leq} \frac{2L_r}{1-\gamma} \mathbb{E}[\|d(\cdot, \cdot; \pi_{\theta_k}) - d(\cdot, \cdot; \pi_{k+1})\|_{TV}] \\
 &\stackrel{(iv)}{\leq} 2L_q C_d \mathbb{E}[\|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|] \\
 &\stackrel{(v)}{\leq} 2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|_{\infty}]
 \end{aligned} \tag{64}$$

where (i) follows the definition $d(s, a; \pi) = (1 - \gamma)\pi(a|s) \sum_{t \geq 0} \gamma^t \mathcal{P}^\pi(s_t = s)$; (ii) is due to distribution mismatch between two visitation measures; (iii) follows the inequality (10a) in Assumption 5.2; the inequality (iv) follows Lemma B.1 and the fact that $\pi_{\theta_k}(\cdot|s) \propto \exp(Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}}(s, \cdot))$, $\pi_{k+1}(\cdot|s) \propto \exp(Q_{r_{\theta_k}, \pi_k}^{\text{soft}}(s, \cdot))$ and the constant $L_q := \frac{L_r}{1-\gamma}$; (v) follows the conversion between Frobenius norm and infinity norm. Through plugging the inequality (64) into (63), it leads to

$$\begin{aligned} & \mathbb{E}[L(\theta_{k+1})] \\ & \geq \mathbb{E}[L(\theta_k)] - 2\alpha L_q \mathbb{E} \left[\left\| \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] - \mathbb{E}_{\tau \sim \pi_{k+1}} \left[\sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta_k) \right] \right\| \right] \\ & \quad + \alpha \mathbb{E} \left[\|\nabla L(\theta_k)\|^2 \right] - 2L_c L_q^2 \alpha^2 \\ & \stackrel{(i)}{\geq} \mathbb{E}[L(\theta_k)] - 4\alpha C_d L_q^2 \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E} \left[\|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|_{\infty} \right] + \alpha \mathbb{E} \left[\|\nabla L(\theta_k)\|^2 \right] - 2L_c L_q^2 \alpha^2 \end{aligned}$$

where (i) follows the inequality (64).

Rearranging the inequality above and denote $C_1 := 4C_d L_q^2 \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}$, it holds that

$$\alpha \mathbb{E} \left[\|\nabla L(\theta_k)\|^2 \right] \leq 2L_c L_q^2 \alpha^2 + \alpha C_1 \mathbb{E} \left[\|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|_{\infty} \right] + \mathbb{E} \left[L(\theta_{k+1}) - L(\theta_k) \right]$$

Summing the inequality above from $k = 0$ to $K - 1$ and dividing both sides by αK , it holds that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla L(\theta_k)\|^2 \right] \leq 2L_c L_q^2 \alpha + \frac{C_1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|Q_{r_{\theta_k}, \pi_{\theta_k}}^{\text{soft}} - Q_{r_{\theta_k}, \pi_k}^{\text{soft}}\|_{\infty} \right] + \mathbb{E} \left[\frac{L(\theta_K) - L(\theta_0)}{K\alpha} \right] \quad (65)$$

Note that the log-likelihood function $L(\theta_K)$ is negative and $L(\theta_0)$ is a bounded constant. Then we could plug (61) into (65), it holds that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla L(\theta_k)\|^2 \right] = \mathcal{O}(K^{-\sigma}) + \mathcal{O}(K^{-1}) + \mathcal{O}(K^{-1+\sigma}) \quad (66)$$

which completes the proof for the inequality (12b). \square