

FEDERATED MAXIMUM LIKELIHOOD INVERSE REINFORCEMENT LEARNING WITH CONVERGENCE GUARANTEE

Anonymous authors

Paper under double-blind review

ABSTRACT

Inverse Reinforcement Learning (IRL) aims to recover the latent reward function and corresponding optimal policy from observed demonstrations. Existing IRL research predominantly focuses on a centralized learning approach, not suitable for real-world problems with distributed data and privacy restrictions. To this end, this paper proposes a novel algorithm for federated maximum-likelihood IRL (F-ML-IRL) and provides a rigorous analysis of its convergence and time-complexity. The proposed F-ML-IRL leverages a dual-aggregation to update the shared global model and performs bi-level local updates – an upper-level learning task to optimize the parameterized reward function by maximizing the discounted likelihood of observing expert trajectories under the current policy and a low-level learning task to find the optimal policy concerning the entropy-regularized discounted cumulative reward under the current reward function. We analyze the convergence and time-complexity of the proposed F-ML-IRL algorithm and show that the global model in F-ML-IRL converges to a stationary point for both the reward and policy parameters within finite time, i.e., the log-distance between the recovered policy and the optimal policy, as well as the gradient of the likelihood objective, converge to zero. Finally, evaluating our F-ML-IRL algorithm on high-dimensional robotic control tasks in MuJoCo, we show that it ensures convergences of the recovered reward in decentralized learning and even outperforms centralized baselines due to its ability to utilize distributed data.

1 INTRODUCTION

Inverse learning is the problem of modeling the preferences and goals of an agent using its observed behavior (Arora & Doshi, 2020). When the behavior of a human expert is observed through demonstration trajectories containing state and action data, Inverse Reinforcement Learning (IRL) models the policy through a Markov Decision Process (MDP) to recover the latent reward function and potentially replicate the human expert’s optimal policy (Russell, 1998). The learned reward function can support various downstream tasks such as agent modeling and transfer learning (Sutton & Barto, 2018; Arora & Doshi, 2020). Recent work has developed provably-efficient IRL algorithms, such as Generative Adversarial Inverse Learning (GAIL) (Ho & Ermon, 2016) and Maximum-likelihood IRL (ML-IRL) (Ratia et al., 2012; Zeng et al., 2022), all using a centralized learning approach. However, demonstration data in practice are often distributed across decentralized clients, e.g., devices, cars, and households. It is not realistic to assume that such sensitive data can always be shared or collected for centralized inverse learning, due to privacy restrictions.

To enable collaborative training of machine learning models among decentralized clients under the privacy restrictions, Federated Learning (FL) provides a promising solution by maintaining training data on local devices and aggregating local updates to build a global model. However, most existing work on FL consider only the forward learning problem, e.g., loss minimization (Li et al., 2019), policy improvement (Jin et al., 2022), learning with heterogeneous models (Zhou et al., 2024), and efficient optimization methods (Li et al., 2020; 2021b;a; Wang et al., 2020), and not the IRL problem. We note that IRL using decentralized clients and distributed data is an open problem. It often has a bi-level structure of maximizing the probability of observing expert trajectories under the current policy and optimizing discounted cumulative reward for the target reward function, which must be solved jointly during IRL. A naive integration of FL and IRL may not achieve convergence. A decentralized learning framework for IRL with theoretical analysis of the convergence and the time-complexity remains a significant challenge.

The goal of this paper is to develop a novel framework for federated maximum-likelihood IRL (F-ML-IRL) and provide a rigorous convergence/time-complexity analysis of the proposed algorithms. We adopt the Maximum Likelihood IRL (ML-IRL) approach in (Zeng et al., 2022) and consider the problem of decentralized IRL of a shared latent reward function, from distributed data and using decentralized client devices. Our solution attains the privacy-preserving benefits of FL in IRL. To address the bi-level nature of IRL, our proposed algorithm’s local training round (McMahan et al., 2017) encompasses an upper-level learning task (on each client with local dataset) to optimize the parameterized reward function to maximize the discounted likelihood of observing expert trajectories under the current policy, as well as a low-level learning task to find the optimal policy concerning the entropy-regularized discounted cumulative reward for the current reward function. Then, we design a dual-aggregation method for aggregating both the action-value networks and reward function models every T local rounds, rather than just one set of model parameters in standard FL. Further, we leverage Soft Q-learning (Haarnoja et al., 2017) as the base RL algorithm. Instead of fully solving the forward RL problem before updating the reward parameter, we perform one-step updates for both the recovered policy and reward parameter alternately to improve the efficiency. To the best of our knowledge, this is the first proposal to formulate and solve this F-ML-IRL problem.

We conduct a rigorous convergence/time-complexity analysis of the proposed F-ML-IRL algorithm. Due to the tight coupling between the reward parameters and the recovered policy in IRL’s bi-level optimization, the dual-aggregation method in our F-ML-IRL must be analyzed to understand its impact on convergence. By bounding the logarithmic distance between the estimated policy and the optimal policy by the distance between their corresponding Q-values, we control the variance introduced by local training by considering the time immediately after each global aggregation. Utilizing the γ -contraction property of soft Q-values, we establish the contraction property of the targeted distance, which allows us to provide a convergence proof for the policy estimate. Moreover, we leverage the Lipschitz continuity of the reward parameter and the convergence of the policy estimate to show that the gradient of the global reward parameter converges to zero as the number of communications increases. These techniques enable us to show that F-ML-IRL’s policy estimation and reward optimization both converge in finite time. The change in convergence speed due to the use of only decentralized clients and distributed data (rather than centralized learning) is characterized.

Our F-ML-IRL is implemented and evaluated on high-dimensional robotic control tasks in MuJoCo (Todorov et al., 2012). We compared its performance with several centralized learning baseline including Behavior Cloning (BC) (Pomerleau, 1988), Generative Adversarial Imitation Learning (GAIL) (Ho & Ermon, 2016), and IRL methods like f-IRL (Ni et al., 2021) and ML-IRL (Zeng et al., 2022). We consider non-iid data distribution, where clients have different local demonstration data with varying performance levels. The baselines are evaluated using centralized data with two setups (i) a single client with medium-level demonstrations and (ii) a single client with a mixture of demonstrations of different levels. The results show that our F-ML-IRL could effectively leverage distributed data and client devices in learning, to achieve similar or better recovered reward than the baselines, while meeting decentralization and data privacy restrictions. Our evaluation code is available at <https://anonymous.4open.science/r/F-ML-IRL/>. The key contributions of this paper are summarized as follows:

- We propose a novel framework for federated maximum-likelihood IRL (F-ML-IRL). It enables decentralized IRL of a shared latent reward function, from distributed data and using decentralized client devices, under data privacy restrictions.
- To support the bi-level optimization structure in IRL – for jointly updating the optimal policy and the reward function estimate, the proposed F-ML-IRL algorithm leverages a dual-aggregation of the model parameters, which ensures convergence to optimal results.
- The convergence and time-complexity of the proposed F-ML-IRL algorithm is quantified, with respect to local rounds T and aggregation steps M . We show that F-ML-IRL achieves convergence in finite time and will have faster convergence with a smaller local rounds T .
- Our solution is evaluated on high-dimensional robotic control tasks in MuJoCo and is shown to achieve similar or higher recovered reward than a number of Imitation Learning and IRL baselines that employ centralized learning.

2 RELATED WORK AND BACKGROUND

IRL aims to learn the reward function using expert demonstration data, which frees the forward RL problem from the requirement of specifying the reward function beforehand (Ng et al., 2000) and also facilitates imitation learning by using the recovered reward function to derive an effective policy

(Abbeel & Ng, 2004). Various formulations and solutions for the IRL problem have been explored. The Maximum Margin Planning algorithm frames the problem within a quadratic programming context Ratliff et al. (2006). Bayesian IRL models infer the posterior distribution of the reward function given a prior (Ramachandran & Amir, 2007). Probabilistic maximum entropy IRL methods favor stochastic policies using entropy regularization. In recent years, Generative Adversarial Imitation Learning (GAIL) (Ho & Ermon, 2016) has adopted a Generative Adversarial Networks (GANs) (Goodfellow et al., 2020) framework to recover the expert’s policy. In this framework, a generator proposes new policies to confuse the discriminator, while the discriminator determines whether the state-action pair from the generator’s policy originates from the expert. However, existing work has not considered the IRL problem with distributed data and decentralized clients, under data privacy.

ML-IRL models the policy through an MDP and recover the latent reward function based on maximum likelihood principle. The convergence of centralized ML-IRL with a single client has been analyzed (Ratia et al., 2012; Zeng et al., 2022) and is shown to outperform other IRL methods. ML-IRL considers a MDP defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \eta, r, \gamma)$, where \mathcal{S} and \mathcal{A} represent the state space and action space, respectively. $\mathcal{P}(s'|s, a)$ denotes the state transition probability, $\eta(\cdot)$ is the initial state distribution, $r(s, a)$ is the reward function, and γ is the discount factor. Let θ denote the parameter vector for the reward function, making the reward function $r(s, a; \theta)$. The IRL problem states that the expert’s behavior is characterized by a stochastic policy $\pi_{r_\theta}(\cdot|s)$. The dataset $\mathcal{D} := \{\tau_m\}_{m=1}^K$ contains trajectories $\tau_m = \{(s_t, a_t)\}_{t=0}^\infty$ from the expert policy $\pi_{r_\theta}(\cdot|s)$.

The discounted log-likelihood of observing all sample trajectories \mathcal{D} from the expert is given by:

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

Assume the state transition probabilities $\mathcal{P}(s_{t+1}|s_t, a_t)$ are known. Then, maximizing the discounted log-likelihood is equivalent to maximizing equation 2.

$$l(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t \geq 0} \gamma^t \log \pi_{r_\theta}(a_t|s_t) \right]. \quad (2)$$

ML-IRL aims to maximizing $l(\theta)$ under the constraint that π_{r_θ} is the optimal policy targeting the discounted cumulative reward regularized by the entropy of the policy, i.e. $\pi_{r_\theta} := \arg \max_{\pi} \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t; \theta) + \mathcal{H}(\pi(\cdot | s_t)))]$, where the entropy of the policy is defined as $\mathcal{H}(\pi(\cdot|s)) := -\sum_{a \in \mathcal{A}} \pi(a|s) \log \pi(a|s)$. Incorporating the policy entropy term as a regularization makes the IRL problem well-defined. This adjustment encourages the agent to explore all possible trajectories in the environment, leading to a more stochastic policy with better generalization capabilities.

For decentralized learning, FL focuses on scenarios where multiple clients work together to train a model using distributed data. FL considers the objective of the form:

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{where} \quad f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w) \quad (3)$$

We assume there are n clients over which the local data D_i is stored. Prior to federated averaging (FedAvg), most works in FL based on Stochastic Gradient Descent (FedSGD) (Shokri & Shmatikov, 2015) ignored the impact of data heterogeneity and imbalance. FedAvg derives from FedSGD but allows multiple rounds of local update $\omega^i \leftarrow \omega^i - \alpha \nabla f_i(\omega^i)$ by gradient descent before aggregating the model parameters at the central server, reducing the frequency and cost of communications. The convergence of FedAvg on non-i.i.d. data has been proved (Li et al., 2019). Since Fed-Avg was proposed as the vanilla FL algorithm, efficient federated optimization methods like FedProx (Li et al., 2020) FedBN (Li et al., 2021b), MOON (Li et al., 2021a), and FedNova (Wang et al., 2020) have been developed to address non-i.i.d. data and accelerate the model training process (Konečný et al., 2016). Additionally, the convergence of model-heterogeneous FL, where reduced-size models are extracted from the global model and applied to low-end clients, was provided in (Zhou et al., 2024). However, existing FL methods could not be directly applied to the ML-IRL problem with decentralized clients, since ML-IRL requires a bi-level optimization involving both policy improvement and reward estimate using maximum likelihood. New algorithms needs to be developed for decentralized ML-IRL with rigorous convergence/time-complexity analysis.

3 FEDERATED MAXIMUM-LIKELIHOOD IRL

3.1 OUR PROBLEM STATEMENT

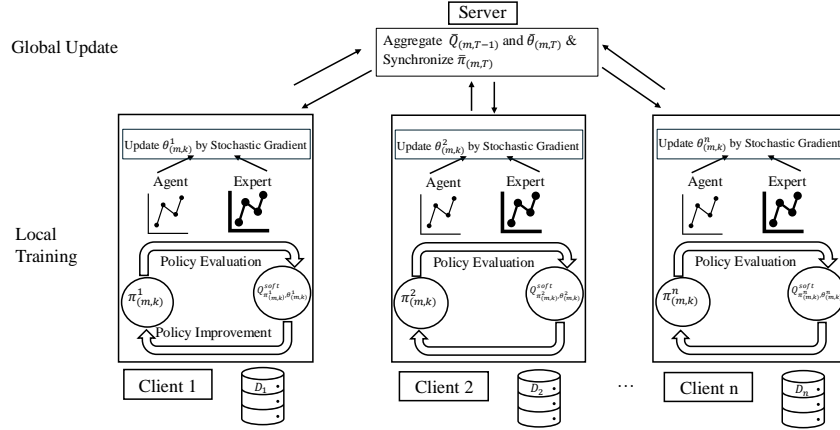


Figure 1: Our F-ML-IRL problem. It aims to recover reward function r_θ from sensitive data/demonstrations $\mathcal{D}_1, \dots, \mathcal{D}_n$ that are distributed over n clients. This requires a novel decentralized algorithm to solve a bi-level optimization – optimizing the parameterized reward function with maximum likelihood and optimizing the corresponding policy concerning the entropy-regularized discounted cumulative reward. We prove the convergence and the time-complexity of F-ML-IRL.

We consider a decentralized inverse learning problem to recover a common reward function r_θ from distributed datasets spread across n clients. Due to privacy requirements, the clients cannot directly share their data for learning. Specifically, we consider n clients, each with a dataset $\mathcal{D}_i := \{\tau_m^i\}_{m=1}^K$ containing trajectories $\tau_m^i = \{(s_t, a_t)\}_{t=0}^\infty$ from the i -th expert policy $\pi_{r_\theta}^i(\cdot|s)$. Different from centralized learning, the clients each have their local model trained on local data. By modeling the distributed expert trajectories as an MDP $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \eta, \gamma)$, our goal is to learn a common reward function r_θ – parameterized by θ – from distributed data and to recover the corresponding optimal policy π_{r_θ} . The F-ML-IRL in this paper is formulated as follows:

$$\begin{aligned} \max_{\theta \in \mathbb{R}^d} \quad & L(\theta) = \frac{1}{n} \sum_{i=1}^n l_i(\theta) \\ \text{s.t.} \quad & \pi_{r_\theta} = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r_\theta(s_t, a_t) + \mathcal{H}(\pi(\cdot | s_t))) \right] \\ \text{where} \quad & l_i(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}_i} \left[\sum_{t \geq 0} \gamma^t \log \pi_{r_\theta}(a_t | s_t) \right] \end{aligned} \quad (4)$$

where $l_i(\theta)$ is the local likelihood calculated using client i 's local data \mathcal{D}_i and target policy π_{r_θ} , which further depends on the current reward function r_θ that is shared by all clients, making it a difficult bi-level optimization. We cannot directly apply FL to this problem, because the maximum likelihood problem on $L(\theta)$ depends on the recovered policy π_{r_θ} , while the policy search for an optimal π_{r_θ} further relies on the estimation of the reward function parameter θ . Thus, the two-level optimization are entangled with each other and requires a new aggregation strategy in F-ML-IRL.

3.2 OUR PROPOSED F-ML-IRL ALGORITHM

We present F-ML-IRL algorithm to solve the decentralized inverse learning problem. Our proposed solution includes three modules - local policy improvement, local reward optimization, and global bi-level aggregation. Each round of F-ML-IRL algorithm consists of T local client steps running in parallel and a global server aggregation of selected model parameters at the end of each round. At each local step, each client i first executes (in parallel) a policy update (on local data \mathcal{D}_i) through policy evaluation and improvement steps based on soft-Q learning to address the lower-level problem. Second, each client carries out a reward optimization, where the reward parameter gradient update is derived by contrasting sampled trajectories from both the expert policy and the current

policy estimate. Next, after every T local steps and at the end of round m , we perform a dual aggregation of both the action-value function and the reward parameters, i.e., to synchronize the local bi-level optimization of both policy and reward on decentralized clients. While our solution is inspired by FL, F-ML-IRL performs a dual aggregation with respect to the bi-level optimization in ML-IRL. The algorithm details are presented below. Its convergence and time-complexity are rigorously analyzed in this paper.

Our F-ML-IRL is illustrated in Fig. 1. Different expert demonstration data D_i are stored at different client devices. We perform local training for policy evaluation and improvement based on soft Q-learning to improve the local policy $\pi_{(m,k)}^i$ under current reward parameter $\theta_{(m,k)}^i$. We then sample trajectories from the current local policy and the expert demonstration data D_i , to provide an update for the reward parameter $\theta_{(m,k)}^i$. At local step k of round m , we use $Q_{\pi_{(m,k)}^i, \theta_{(m,k)}^i}^{\text{soft}}(s, a)$ to denote the action-value function (i.e., Q-value) for action a and state s , with respect to the current policy estimation $\pi_{(m,k)}^i$ under current reward parameter estimation $\theta_{(m,k)}^i$, on each client i . After every T steps of local training, we perform dual aggregation of the Q-values $\bar{Q}_{(m,T-1)}^{\text{soft}}$ and the reward parameter $\bar{\pi}_{(m,T)}$. To the best of our knowledge, this is the first paper considering an ML-IRL problem in this FL context.

Local training for policy improvement. Iterations of local training on each local client start with a shared model with parameters $\pi_{(m,0)}^i(\cdot|s)$ and $\theta_{(m,0)}^i$. During each local training round, we first evaluate the local policy $\pi_{(m,k)}^i(\cdot|s)$ by computing the Q-values $Q_{(m,k)}^i(\cdot, \cdot)$ under the fixed reward parameter θ^i for the i -th local client using the definitions of the soft value and soft Q functions in equation 5.

$$\begin{aligned} V_{\pi_{(m,k)}^i, \theta_{(m,k)}^i}^{\text{soft}}(s) &= \mathbb{E}_{\pi_{(m,k)}^i, s_0=s} \sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t; \theta_{(m,k)}^i) + \mathcal{H}(\pi_{\theta_{(m,k)}^i}(\cdot|s_t)) \right) \\ Q_{\pi_{(m,k)}^i, \theta_{(m,k)}^i}^{\text{soft}}(s, a) &= r(s, a; \theta_{(m,k)}^i) + \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s,a)} V_{\pi_{(m,k)}^i, \theta_{(m,k)}^i}^{\text{soft}}(s') \end{aligned} \quad (5)$$

Then, $\pi_{(m,k+1)}^i(\cdot|s)$ is updated according to the policy improvement step using soft Q-learning in equation 6. It does not assume an explicit policy function but uses the Boltzmann distribution of the Q function, making the probability of choosing an action at some state s proportional to the exponential of the Q-value of this action-state pair.

$$\pi_{(m,k+1)}^i(a|s) \propto \exp(Q_{\pi_{(m,k)}^i, \theta_{(m,k)}^i}^{\text{soft}}(s, a)), \forall s \quad (6)$$

Local training for reward optimization. For the optimization towards the local reward parameter $\theta_{(m,k+1)}^i$, a stochastic gradient ascent method is proposed. The gradient of each local likelihood function $l_i(\theta)$ is given by equation 7, which derives from Lemma 1 in (Zeng et al., 2022).

$$\nabla l_i(\theta) = \mathbb{E}_{\tau_i \sim \mathcal{D}_i} \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta) - \mathbb{E}_{\tau_i \sim \pi_{\theta}} \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta). \quad (7)$$

We construct a stochastic estimator of the exact gradient $\nabla l_i(\theta_{(m,k)}^i)$, approximating the optimal policy $\pi_{r_{\theta_{(m,k)}^i}}$ with the current policy $\pi_{(m,k+1)}^i$. Specifically, we sample one expert trajectory $\tau_{(m,k)}^{E_i} := \{s_t, a_t\}_{t \geq 0}$ from the local dataset \mathcal{D}_i and one agent trajectory $\tau_{(m,k)}^{A_i} := \{s_t, a_t\}_{t \geq 0}$ from the current policy $\pi_{(m,k+1)}^i$. Then we use a stochastic estimate $g_{(m,k)}^i$ to approximate the exact gradient of the local likelihood objective function l_i for each local client in equation 8. The update of the reward relies on both the local softmax policy $\pi_{(m,k+1)}^i$ through $\tau_{(m,k)}^{A_i}$ and the local data \mathcal{D}_i through $\tau_{(m,k)}^{E_i}$.

$$g_{(m,k)}^i = h(\theta_{(m,k)}^i; \tau_{(m,k)}^{E_i}) - h(\theta_{(m,k)}^i; \tau_{(m,k)}^{A_i}) \quad (8)$$

where $h(\theta; \tau) = \sum_{t \geq 0} \gamma^t \nabla_{\theta} r(s_t, a_t; \theta)$. Finally, the local reward parameter $\theta_{(m,k)}^i$ is updated as:

$$\theta_{(m,k+1)}^i = \theta_{(m,k)}^i + \alpha g_{(m,k)}^i \quad (9)$$

where α is the learning rate of the reward parameter update.

Bi-level model aggregation. Every T local iterations, local Q-values and local reward parameters are communicated to the global server for aggregation, while the policy synchronization is performed based on the aggregated Q-values such that each local client has the same policy after the aggregation. We design the dual aggregation step after thorough thoughts. The reward update in equation 9 depends on how well the trajectories from policy $\pi_{(m,k)}^i$ approximates the optimal policy $\pi_{r_{\theta_{(m,k)}^i}}$, while the policy $\pi_{(m,k)}^i$ relies on the Q-value update from equation 5. Therefore, our FL algorithm aims to improve the Q-value estimates for local clients by aggregating their Q-values equation 10.

$$\bar{Q}_{(m,T-1)}^{\text{soft}}(\cdot, \cdot) := \sum_{j=1}^N Q_{(m,T-2)}^j(\cdot, \cdot) / N \quad (10)$$

We note that when the Q-values are represented by another network with parameter ψ , the aggregation of the Q-values will simply become aggregation of model parameters. The policy synchronization is automatically performed by policy improvement based on the aggregated Q-values and sent to each local client for update such that each local client has the same policy after the Q aggregation in equation 11:

$$\bar{\pi}_{(m,T)}(\cdot | s) \propto \exp(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \cdot)), \forall s \in S \quad (11)$$

Since ML-IRL requires a bi-level problem with respect to both the reward parameter and the recovered policy, we consider a dual aggregation that also applies to the reward parameter θ :

$$\bar{\theta}_{(m,T)} := \sum_{j=1}^N \theta_{(m,T-1)}^j / N \quad (12)$$

After each dual aggregation, the global policy and reward parameters are sent to each local clients as an initialization for future local training: $\pi_{(m,0)}^i(\cdot | s) = \bar{\pi}_{(m-1,T)}(\cdot | s)$ and $\theta_{(m,0)}^i = \bar{\theta}_{(m-1,T)}$ for all $i = 1, 2, \dots, N$. The entire process of the F-ML-IRL algorithm is summarized in Algorithm 1.

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

- 1: **Input:** Initialize reward parameter $\theta_{(0,0)}^i$ and policy $\pi_{(0,0)}^i$. Set the aggregation period to be T , number of local server to be N , and reward parameter's local stepsize as α .
 - 2: **for** $m = 0, 1, \dots, M - 1$ **do**
 - 3: **if** $m > 0$ **then**
 - 4: Inherit $\pi_{(m,0)}^i(\cdot | s)$ and $\theta_{(m,0)}^i$ from last aggregation
 - 5: **end if**
 - 6: **for** $k = 0, \dots, T - 2$ **do**
 - 7: **for** $i = 1, 2, \dots, N$ **do**
 - 8: Compute $Q_{r_{\theta_{(m,k)}^i}, \pi_{(m,k)}^i}^{\text{soft}}(\cdot, \cdot)$ using equation 5
 - 9: Update $\pi_{(m,k+1)}^i(\cdot | s)$ based on equation 6
 - 10: Sample an expert trajectory $\tau_{(m,k)}^{E_i}$ from local dataset D_i
 - 11: Sample a trajectory $\tau_{(m,k)}^{A_i}$ from current policy $\pi_{(m,k+1)}^i$
 - 12: Estimating gradient $g_{(m,k)}^i$ following equation 8
 - 13: Update reward parameter $\theta_{(m,k+1)}^i$ using equation 9
 - 14: **end for**
 - 15: **end for**
 - 16: **Set** $k = T - 1$
 - 17: Aggregate $\bar{Q}_{(m,k)}^{\text{soft}}(\cdot, \cdot)$ by equation 10
 - 18: Synchronize policies $\bar{\pi}_{(m,k+1)}(\cdot | s)$ using equation 11
 - 19: Aggregate reward parameters $\bar{\theta}_{(m,k+1)}$ by equation 12
 - 20: **end for**
-

4 THEORETICAL ANALYSIS

4.1 ASSUMPTIONS

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

$$\sup_{s \in \mathcal{S}} \|\mathbb{P}(s_t \in \cdot \mid s_0 = s, \pi) - \mu_\pi(\cdot)\|_{TV} \leq \kappa \rho^t, \quad \forall t \geq 0 \quad (13)$$

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

Equation 13 states that the Markov chain mixes at a geometric rate. This is a common assumption in the RL literature, which holds for any time-homogeneous Markov chain with a finite state space.

Bounded Gradient and Lipschitz Property. For any $s \in \mathcal{S}$, $a \in \mathcal{A}$, and any reward parameter θ , the following conditions hold, where L_r and L_g are positive constants:

$$\|\nabla_{\theta} r(s, a; \theta)\| \leq L_r, \quad \text{and} \quad \|\nabla_{\theta} r(s, a; \theta_1) - \nabla_{\theta} r(s, a; \theta_2)\| \leq L_g \|\theta_1 - \theta_2\|, \quad (14)$$

Equation 14 posits that the parameterized reward function has a bounded gradient and is Lipschitz smooth, which is common in the literature.

4.2 IMPORTANT LEMMAS

We first introduce two important lemmas that are used repeatedly in the converge analysis. Due to space limitations, the proofs of these lemmas are included in the appendix.

Lemma 1. *Suppose the above assumptions hold. Given any reward parameters θ_1 and θ_2 , the following results hold for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$:*

$$\left| Q_{r_{\theta_1}, \pi_{\theta_1}}^{\text{soft}}(s, a) - Q_{r_{\theta_2}, \pi_{\theta_2}}^{\text{soft}}(s, a) \right| \leq L_q \|\theta_1 - \theta_2\|, \quad (15)$$

$$\|\nabla L(\theta_1) - \nabla L(\theta_2)\| \leq L_c \|\theta_1 - \theta_2\|, \quad (16)$$

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} .

Lemma 1 is directly derived from Lemma 2 in (Zeng et al., 2022), where the positive constants L_q and L_c are also defined. The Lipschitz properties of the Q-value function and the gradient of the log-likelihood are essential for convergence analysis, as they help control the distance between local and global models in the FL setting.

Lemma 2. *For any two policies $\pi(a|s)$ and $\pi'(a|s)$, the difference in their soft Q-values under some reward function r for a given state-action pair (s, a) is bounded as follows:*

$$\|Q_{r_{\theta}, \pi}^{\text{soft}} - Q_{r_{\theta}, \pi'}^{\text{soft}}\|_{\infty} \leq \frac{\gamma}{1 - \gamma} \|\log(\pi) - \log(\pi')\|_{\infty} \quad (17)$$

Controlling the distance between soft Q-values under different policies helps us analyze the optimality of the global policy with respect to the global reward parameter after aggregations.

4.3 MAIN CONVERGENCE RESULT

Theorem 1. *Under the above two assumptions, if we choose step size $\alpha_{(m,k)} = \alpha_0 / (mT + k)^{\sigma}$ in F-ML-IRL (Algorithm 1), where $\alpha_0 > 0$ and $\sigma \in (0, 1)$ are constants and M is the total number of dual aggregations, the following convergence results hold for F-ML-IRL:*

$$\frac{1}{M} \sum_{m=0}^{M-1} \mathbb{E} \left[\left\| \log(\bar{\pi}_{(m,T)}) - \log(\pi_{\theta^i_{(m,T-1)}}) \right\|_{\infty} \right] = \mathcal{O}(M^{-1} \gamma^{T-1}) + \mathcal{O}(M^{-\sigma} T^{1-\sigma}) \quad (18)$$

$$\frac{1}{M} \sum_{m=0}^{M-1} \mathbb{E} \left[\left\| \nabla L(\bar{\theta}_{(m,T)}) \right\|^2 \right] = \mathcal{O}(M^{-1}) + \mathcal{O}(M^{-\sigma} T^{-\sigma}) + \mathcal{O}(M^{-1-\sigma} T^{1-\sigma}). \quad (19)$$

Remarks: The time complexity of both policy estimate and reward parameter optimization depends on the number of global aggregation rounds M and the number of local training steps T . The policy and reward function parameters in F-ML-IRL converge at the rate of $M^{-\sigma}T^{1-\sigma}$ and $M^{-\sigma}T^{-\sigma}$, respectively, since we have $\sigma \in (0, 1)$ and T is often fixed. We note that due to dual-aggregation and the variance caused by local training on distributed datasets across decentralized clients, F-ML-IRL exhibits a slightly slower convergence rate, compared with standard centralized ML-IRL with a single client (whose convergence rate is $M^{-\sigma}$). From Equations (18) and (19), there exists a sweet spot with respect to the number of local training steps T , since γ^{T-1} and $T^{-\sigma}$ both decreases with T , while $T^{1-\sigma}$ increases. Exploring this tradeoff will be considered in future work. Compared with Fed-Avg (whose convergence rate is $M^{-1}T^{-1}$), F-ML-IRL also has a slower convergence rate due to the complexity of the bi-level optimization problem.

Proof Sketch: Due to space limitations, we outline the key steps of our convergence analysis and present the complete proof in the appendix. We first analyze the convergence of policy estimates and reduce it to the convergence of Q-values. We then analyze the distance between Q-values using the Lipschitz property, tracing back to the start of each dual aggregation around. In particular, we examine the extra distance between the estimated policy and the optimal policy caused by aggregation, seeking the contraction property of Q-value estimates between adjacent aggregation rounds. Next, for reward optimization, we leverage the Lipschitz smooth property of the likelihood and control the discrepancy between the stochastic gradient and the true gradient. This allows us to use the convergence of Q-values from the previous analysis to demonstrate the gradient convergence of the reward parameter. For simplicity of notations, we use $Q_{i,(m,t)}^{\text{soft}}$ to denote $Q_{r_{\theta^i(m,t)}, \pi^i(m,t)}$, the action-value function at a given state for the local policy and reward parameter estimations at round (m, t) . Similarly, $Q_{i,(m,t)}^{\text{soft}*}$ denotes $Q_{r_{\theta^i(m,t)}, \pi^i(m,t)}$, which is the Q-function for the optimal policy under the reward parameter at round (m, t) and $Q_{(m,t)}^{\text{soft}}$ denotes $Q_{r_{\bar{\theta}^i(m,T)}, \pi_{\bar{\theta}^i(m,T)}}$, which represents the Q-function for the aggregated policy and reward parameter at the m 'th aggregation.

Convergence of Policy Estimate:

Step 1: We show the distance between the aggregated policy and the optimal policy under the pre-aggregation local reward parameter could be controlled using the distance between soft Q-functions:

$$\|\log(\bar{\pi}_{(m,T)}) - \log(\pi_{\theta^i(m,T-1)}^i)\|_{\infty} \leq 2\|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{i,(m,T-2)}^{\text{soft}}\|_{\infty}$$

This step relies on the policy update rule in equation 6.

Step 2: By introducing intermediary terms as bridges, specifically looking back to the time right after the last aggregation, where all local servers have the same reward parameter $\bar{\theta}_{(m-1,T)}$, we further bound the difference in **Step 1** by converting it to the difference of reward parameters using equation 9, 10, 15. Combining this with the γ -contraction property of the soft-Q update, we have:

$$\|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{i,(m,T-2)}^{\text{soft}}\|_{\infty} \leq \gamma^{T-2}\|Q_{r_{\bar{\theta}_{(m-1,T)}}, \bar{\pi}_{(m-1,T)}}^{\text{soft}} - Q_{r_{\bar{\theta}_{(m-1,T)}}, \pi_{\bar{\theta}_{(m-1,T)}}^i}\|_{\infty} + E_1$$

where we use auxiliary variable $E_1 = 4\alpha \left(\frac{1-\gamma^{T-2}}{1-\gamma} + T - 2 \right) L_q^2$.

Step 3: Using Lemma 2, we bound the difference in Q-values corresponding to different policies with the same reward during the aggregation step, and finally have:

$$\|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{i,(m,T-2)}^{\text{soft}}\|_{\infty} \leq (1-\gamma)\gamma^{T-2}\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{i,(m-1,T-2)}^{\text{soft}}\|_{\infty} + E_2$$

where we use auxiliary variable $E_2 = 2\frac{1-\gamma^{T-2}}{1-\gamma} + (1-\gamma)^2\gamma^{T-2} + \frac{1-\gamma}{\gamma}(2T-3) + 2(T-2)L_q^2$. Finally, we obtain the convergence rate of the policy estimate by the contraction of Q-difference.

Convergence of Reward Parameter Optimization:

Step 1: We first leverage the Lipschitz smooth property of $l(\theta)$ equation 16:

$$L(\bar{\theta}_{(m,T)}) \geq L(\bar{\theta}_{(m-1,T)}) + \langle \nabla L(\bar{\theta}_{(m,T)}), \bar{\theta}_{(m,T)} - \bar{\theta}_{(m-1,T)} \rangle - \frac{L_c}{2}\|\bar{\theta}_{(m,T)} - \bar{\theta}_{(m-1,T)}\|^2$$

Step 2: We show the bias between the stochastic gradient estimate $g_{(m,k)}^i$ and the true gradient $\nabla L(\theta_{(m-1,T)})$ could be controlled. In this process, we also compare the increments of local clients

to control the extra error terms introduced by the federated scheme leveraging equation 9, 12. We show that the gradient of the global reward parameter could be bounded by the distance between Q-values:

$$\alpha(T-1)\mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] \leq \alpha C_1 \mathbb{E}[\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{i,(m-1,T-2)}^{\text{soft}}\|_\infty] + \mathbb{E}[L(\bar{\theta}_{(m,T)}) - L(\bar{\theta}_{(m-1,T)})] + E_3 \quad (20)$$

where $C_1 = \frac{4(1-\gamma^{T-1})}{\gamma} L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}$ and $E_3 = 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1-\frac{1-\gamma^{T-1}}{1-\gamma}}{1-\gamma} + \frac{(T-1)(3T-1)\alpha^2 L_c L_q^2}{2} + \frac{4(1-\gamma^{T-1})}{1-\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot [2(2T-3)\alpha L_q^2 + \frac{1-\gamma}{\gamma} \cdot 4\alpha L_q^2]$ are two auxiliary variables. By combining this with the convergence of the Q-value difference that was established in **Step 3** of the Policy Estimation proof, we obtain the desired convergence of the reward parameter.

5 EVALUATIONS

We evaluated the proposed F-ML-IRL method on five high-dimensional robotic control tasks in MuJoCo (Todorov et al., 2012). For comparison, we selected several state-of-the-art baselines, including imitation learning approaches that only learn the expert policy—specifically like BC (Pomerleau, 1988) and GAIL (Ho & Ermon, 2016), as well as IRL methods that simultaneously learn both a reward function and a policy, namely f-IRL (Ni et al., 2021) and ML-IRL (Zeng et al., 2022). To ensure fairness, we used Soft Actor-Critic (SAC) (Haarnoja et al., 2018) as the base RL algorithm for all methods since it incorporates elements of Soft Q-Learning and achieves strong performance using the actor-critic scheme. The experiments are conducted on a server with AMD EPYC 7513 32-Core Processors and NVIDIA RTX A6000 GPUs. We choose $M = 200$ rounds and $T = 5$ local steps and average the results over multiple runs. Our evaluation code is available at <https://anonymous.4open.science/r/F-ML-IRL/>.

Environment	Setting	F-ML-IRL	ML-IRL		BC		GAIL		f-IRL	
			Mixed	Medium	Mixed	Medium	Mixed	Medium	Mixed	Medium
Ant	(3, 200)	6425.91	6219.78	6161.65	983.99	984.04	989.30	988.73	5615.33	5930.89
	(3, 1000)	6398.98	5100.25	6402.87	5952.08	718.87	989.00	989.29	5370.28	5527.63
	(5, 200)	6254.32	5614.91	6161.65	983.51	984.04	988.67	988.73	5628.94	5930.89
	(5, 1000)	6528.04	6330.67	6402.87	411.83	718.87	989.77	989.29	5388.74	5527.63
HalfCheetah	(3, 200)	13007.75	8054.94	12581.28	-0.63	-0.73	7513.31	10288.42	10110.73	12962.52
	(3, 1000)	13228.98	13642.82	13124.24	-0.66	-11.74	12112.99	11506.59	13075.95	12871.64
	(5, 200)	11827.91	6406.45	12581.28	-0.57	-0.73	4910.45	10288.42	7132.01	12962.52
	(5, 1000)	12360.60	12750.04	13124.24	110.59	-11.74	11364.40	11506.59	12659.30	12871.64
Hopper	(3, 200)	3576.10	1871.83	3623.07	18.11	18.13	1022.90	1023.65	1297.25	3456.47
	(3, 1000)	3674.64	3518.04	3479.33	18.17	2290.58	1025.93	1032.07	3403.36	3390.08
	(5, 200)	3419.95	1484.52	3623.07	18.12	18.13	1020.19	1023.65	1313.12	3456.47
	(5, 1000)	3618.44	3601.40	3479.33	1016.31	2290.58	1111.08	1032.07	3468.72	3390.08
Humanoid	(3, 200)	5656.06	5484.99	5861.01	243.21	242.50	4666.38	3035.34	5510.06	6004.58
	(3, 1000)	5694.79	5903.42	5813.57	241.46	532.41	4627.15	4688.41	5708.21	5726.86
	(5, 200)	6232.37	5462.64	5861.01	242.69	242.50	4692.86	3035.34	5523.40	6004.58
	(5, 1000)	6294.25	5713.37	5813.57	545.01	532.41	4577.12	4688.41	5608.83	5726.86
Walker2d	(3, 200)	4057.25	3317.61	4400.43	8.38	8.27	353.66	18.74	1050.78	5729.55
	(3, 1000)	5798.37	5061.23	5673.49	8.27	507.69	344.55	19.27	4805.53	5255.57
	(5, 200)	4540.90	3024.14	4400.43	8.40	8.27	13.03	18.74	1115.37	5729.55
	(5, 1000)	5853.42	4669.73	5673.49	711.06	507.69	360.10	19.27	4704.77	5255.57
Average	-	6712.60	5661.64	6712.09	575.97	529.00	3183.64	3359.05	5424.58	6685.58

Table 1: Compare F-ML-IRL and baselines on MuJoCo tasks, with different number of clients and demonstration trajectory length. F-ML-IRL achieves similar or higher recovered reward in almost all scenarios and outperforms the baselines in more than half, as well as in terms of the average.

We evaluate different algorithms using the rewards associated with the recovered expert policies evaluated in the original environment (same as the method adopted in previous work). We compare F-ML-IRL with the baselines on five MuJoCo tasks under non-iid data distributions, where each

client contains different demonstration data corresponding to varying levels of expertise. For the baselines that rely on centralized learning, we consider two setups: (i) a single client with medium-level demonstrations, denoted as *medium* and (ii) a single client with a mixture of demonstrations of different levels, denoted as *mixed*. In either case, the total amount of local data per client remains the same in the experiments. More details on experiment set up is provided in the appendix.

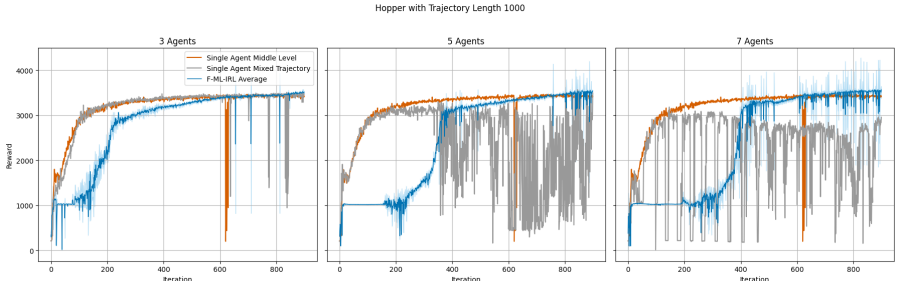


Figure 2: Convergence of F-ML-IRL in Hopper Environment compared with centralized ML-IRL with mixed and medium data. As the number of clients (and thus the non-iid datasets) increases from 3 (left) to 7 (right), F-ML-IRL takes longer to converge and nevertheless achieves more significant improvement by leveraging distributed demonstration data on the clients.

The evaluation results are summarized in Table 1. We have tested each algorithm and each MuJoCo task under 4 settings, i.e., with 3 or 5 clients and with demonstration trajectory length equals to 200 or 1000, respectively. As demonstrated in (Zeng et al., 2022), even a single expert trajectory of length 1000 can lead to a well-recovered policy using ML-IRL. To investigate the performance of our model under conditions of scarce and distributed data, we utilize a single expert trajectory of length 1000 and further reduce its length to 200 in the experiments. In Table 1, we also compute the average reward for each algorithm across all settings and tasks in our experiments.

We note that F-ML-IRL ensures convergences of the recovered reward in decentralized learning and achieves similar or higher recovered reward than the baselines in almost all settings and tasks. It even outperforms centralized baselines in more than half of the settings and tasks, due to its ability to utilize distributed data. The performance of F-ML-IRL is pretty robust as the number of clients increases to 5 and the expert trajectory length reduces to 200. Imitation learning baselines like BC and GAIL generally have lower performance and even fail in some settings. While ML-IRL performs generally well, it fails to recover a satisfactory policy when data is limited or in tasks involving mixed trajectories of different expertise. On the other hand, f-IRL performs relatively well when provided with longer expert trajectories but struggles when demonstration data is limited. In contrast, our F-ML-IRL consistently achieves similar or higher recovered rewards compared to all baselines, particularly maintaining robust performance even when data is limited and involves demonstrations of mixed expertise.

We further illustrate the convergence of our F-ML-IRL algorithm compared with two different centralized learning baselines using ML-IRL (with medium and mixed-data, respectively) in the Hopper environment, as shown in Figure 2. As the number of clients (and thus the number of non-iid local datasets) increases (from 3 clients on the left to 7 clients on the right), it takes F-ML-IRL more rounds to converge, because of the increased variance introduced by local training on more participating clients and datasets. Nevertheless, F-ML-IRL is able to converge to higher recovered reward than both baselines. Centralized ML-IRL suffers with mixed demonstration data of varying expertise. In contrast, as the number of clients and demonstration dataset increases, F-ML-IRL shows more significant improvement by leverage distributed demonstration data on the clients.

6 CONCLUSIONS

This paper proposes F-ML-IRL for federated maximum-likelihood inverse reinforcement learning. It enables decentralized learning of a shared latent reward function from distributed datasets and using decentralized clients. F-ML-IRL algorithm leverages a dual-aggregation to update the shared global model and performs bi-level local updates for inverse learning. We analyze the convergence and time-complexity of F-ML-IRL. Evaluation results on MuJoCo tasks how that F-ML-IRL ensures convergences of the recovered reward and achieves similar or higher recovered reward, compared to state-of-the-art baselines using centralized inverse learning. For further work, we plan to investigate further communication reduction and the use of heterogeneous local models in F-ML-IRL.

REFERENCES

- 540
541
542 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In
543 *Proceedings of the twenty-first international conference on Machine learning*, pp. 1, 2004.
- 544 Saurabh Arora and Prashant Doshi. A survey of inverse reinforcement learning: Challenges, meth-
545 ods and progress, 2020. URL <https://arxiv.org/abs/1806.06877>.
- 546 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
547 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the*
548 *ACM*, 63(11):139–144, 2020.
- 549 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with
550 deep energy-based policies. In *International conference on machine learning*, pp. 1352–1361.
551 PMLR, 2017.
- 552 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
553 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
554 *ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 555 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
556 *information processing systems*, 29, 2016.
- 557 Hao Jin, Yang Peng, Wenhao Yang, Shusen Wang, and Zhihua Zhang. Federated reinforcement
558 learning with environment heterogeneity. In *International Conference on Artificial Intelligence*
559 *and Statistics*, pp. 18–37. PMLR, 2022.
- 560 Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization:
561 Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016.
- 562 Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of*
563 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10713–10722, 2021a.
- 564 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.
565 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and sys-*
566 *tems*, 2:429–450, 2020.
- 567 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of
568 fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- 569 Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning
570 on non-iid features via local batch normalization. *arXiv preprint arXiv:2102.07623*, 2021b.
- 571 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
572 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*
573 *gence and statistics*, pp. 1273–1282. PMLR, 2017.
- 574 Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*, vol-
575 ume 1, pp. 2, 2000.
- 576 Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Ben Eysenbach. f-irl: Inverse
577 reinforcement learning via state marginal matching. In *Conference on Robot Learning*, pp. 529–
578 551. PMLR, 2021.
- 579 Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in neural*
580 *information processing systems*, 1, 1988.
- 581 Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *IJCAI*, vol-
582 ume 7, pp. 2586–2591, 2007.
- 583 Héctor Ratia, Luis Montesano, and Ruben Martinez-Cantin. On the performance of maximum
584 likelihood inverse reinforcement learning. *arXiv preprint arXiv:1202.1558*, 2012.
- 585 Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. Maximum margin planning. In
586 *Proceedings of the 23rd international conference on Machine learning*, pp. 729–736, 2006.
- 587 Stuart Russell. Learning agents for uncertain environments. In *Proceedings of the eleventh annual*
588 *conference on Computational learning theory*, pp. 101–103, 1998.

594 Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd*
595 *ACM SIGSAC conference on computer and communications security*, pp. 1310–1321, 2015.
596
597 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
598 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
599 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033.
600 IEEE, 2012.
601 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective
602 inconsistency problem in heterogeneous federated optimization. *Advances in neural information*
603 *processing systems*, 33:7611–7623, 2020.
604 Siliang Zeng, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Maximum-likelihood inverse re-
605 inforcement learning with finite-time guarantees. *Advances in Neural Information Processing*
606 *Systems*, 35:10122–10135, 2022.
607 Hanhan Zhou, Tian Lan, Guru Prasad Venkataramani, and Wenbo Ding. Every parameter matters:
608 Ensuring the convergence of federated learning with dynamic heterogeneous models reduction.
609 *Advances in Neural Information Processing Systems*, 36, 2024.
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 A PROOF OF LEMMA 2

649 Given the definition of soft-Q function following Bellman equation:

$$650 \quad Q_{r,\pi}^{\text{soft}}(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} [r(s, a) + \gamma \mathbb{E}_{a' \sim \pi(\cdot|s')} (Q_{r,\pi}^{\text{soft}}(s', a') - \log \pi(a'|s'))] \quad (21)$$

651 The difference between soft-Q values under policies π and π' is:

$$652 \quad |Q_{r,\pi}^{\text{soft}}(s, a) - Q_{r,\pi'}^{\text{soft}}(s, a)| = \gamma \left| \mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathbb{E}_{a' \sim \pi(\cdot|s')} (Q_{r,\pi}^{\text{soft}}(s', a') - \log \pi(a'|s'))] \right. \\ 653 \quad \left. - \mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathbb{E}_{a' \sim \pi'(\cdot|s')} (Q_{r,\pi'}^{\text{soft}}(s', a') - \log \pi'(a'|s'))] \right| \quad (22)$$

654 Using the triangle inequality, we separate the terms in equation 22:

$$655 \quad |Q_{r,\pi}^{\text{soft}}(s, a) - Q_{r,\pi'}^{\text{soft}}(s, a)| \leq \gamma \left(\left| \mathbb{E}_{s' \sim P(\cdot|s,a)} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q_{r,\pi}^{\text{soft}}(s', a') - Q_{r,\pi'}^{\text{soft}}(s', a')] \right| \right. \\ 656 \quad \left. + \left| \mathbb{E}_{s' \sim P(\cdot|s,a)} \mathbb{E}_{a' \sim \pi'(\cdot|s')} [\log \pi(a'|s') - \log \pi'(a'|s')] \right| \right) \quad (23)$$

657 For the first term in equation 22, we apply Jensen's inequality to the absolute value function:

$$658 \quad |\mathbb{E}_{s' \sim P(\cdot|s,a)} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q_{r,\pi}^{\text{soft}}(s', a') - Q_{r,\pi'}^{\text{soft}}(s', a')]| \leq \sup_{s', a'} |Q_{r,\pi}^{\text{soft}}(s', a') - Q_{r,\pi'}^{\text{soft}}(s', a')| \quad (24)$$

659 Similarly, the second term in equation 22 involving the log policies is bounded as:

$$660 \quad |\mathbb{E}_{s' \sim P(\cdot|s,a)} \mathbb{E}_{a' \sim \pi'(\cdot|s')} [\log \pi(a'|s') - \log \pi'(a'|s')]| \leq \sup_{s', a'} |\log \pi(a'|s') - \log \pi'(a'|s')| \quad (25)$$

661 Thus, we combine the two terms in equation 22:

$$662 \quad |Q_{r,\pi}^{\text{soft}}(s, a) - Q_{r,\pi'}^{\text{soft}}(s, a)| \leq \gamma \sup_{s', a'} |Q_{r,\pi}^{\text{soft}}(s', a') - Q_{r,\pi'}^{\text{soft}}(s', a')| + \gamma \sup_{s', a'} |\log \pi(a'|s') - \log \pi'(a'|s')| \quad (26)$$

663 Since soft-Q values depend recursively on future rewards, we apply this bound recursively over time. At $t = 1$, the same bound holds:

$$664 \quad |Q_{r,\pi}^{\text{soft}}(s_1, a_1) - Q_{r,\pi'}^{\text{soft}}(s_1, a_1)| \leq \gamma \sup_{s_2, a_2} |Q_{r,\pi}^{\text{soft}}(s_2, a_2) - Q_{r,\pi'}^{\text{soft}}(s_2, a_2)| \\ 665 \quad + \gamma \sup_{s', a'} |\log \pi(a'|s') - \log \pi'(a'|s')| \quad (27)$$

666 Substituting this into the previous equation, we get:

$$667 \quad |Q_{r,\pi}^{\text{soft}}(s_0 = s, a_0 = a) - Q_{r,\pi'}^{\text{soft}}(s_0 = s, a_0 = a)| \\ 668 \quad \leq \gamma^2 \sup_{s_2, a_2} |Q_{r,\pi}^{\text{soft}}(s_2, a_2) - Q_{r,\pi'}^{\text{soft}}(s_2, a_2)| + \gamma \sup_{s', a'} |\log \pi(a'|s') - \log \pi'(a'|s')| (1 + \gamma) \quad (28)$$

702 Applying this recursively over n steps:

$$703 \quad |Q_{r,\pi}^{\text{soft}}(s_0 = s, a_0 = a) - Q_{r,\pi'}^{\text{soft}}(s_0 = s, a_0 = a)|$$

$$704 \quad \leq \gamma^n \sup_{s_n, a_n} |Q_{r,\pi}^{\text{soft}}(s_n, a_n) - Q_{r,\pi'}^{\text{soft}}(s_n, a_n)| + \sup_{s', a'} |\log \pi(a'|s') - \log \pi'(a'|s')| \sum_{k=0}^{n-1} \gamma^k \quad (29)$$

709 As $n \rightarrow \infty$, the term $\gamma^n \sup_{s_n, a_n} |Q_{r,\pi}^{\text{soft}}(s_n, a_n) - Q_{r,\pi'}^{\text{soft}}(s_n, a_n)|$ tends to zero. The geometric series
710 sum is $\sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$. Thus, taking the infinity norm with respect to all s and a , the final bound
711 is:

$$712 \quad \|Q_{r,\pi}^{\text{soft}} - Q_{r,\pi'}^{\text{soft}}\|_{\infty} \leq \frac{\gamma}{1-\gamma} \|\log(\pi) - \log(\pi')\|_{\infty} \quad (30)$$

717 B PROOF OF THEOREM 1

718 In this section we show the complete proof for the convergence results for the recovered global poli-
719 cies $\bar{\pi}_{(m,T)}$ and the global reward parameter $\bar{\theta}_{(m,T)}$ after m communications with communication
720 period T .

722 B.1 CONVERGENCE OF POLICY ESTIMATE $\bar{\pi}_{(m,T)}$

723 We first analyze the approximation error between the logarithm of the synchronized policy
724 $\log(\bar{\pi}_{(m,T)}(a|s))$ and the logarithm of the optimal policy corresponding to the previous local re-
725 ward parameter $\log(\pi_{\theta_{(m,T-1)}^i}(a|s))$ for all i . Specifically, we aim to bound the difference:

$$726 \quad \left| \log(\bar{\pi}_{(m,T)}(a|s)) - \log(\pi_{\theta_{(m,T-1)}^i}(a|s)) \right|$$

727 This difference represents the discrepancy between the synchronized policy after the m -th global
728 aggregation and the optimal policy corresponding to the previous local reward parameter $\theta_{(m,T-1)}^i$.

729 We aim to show that the distance between the logarithms of the synchronized policy and the optimal
730 policy can be bounded by the difference between their corresponding soft-Q values. Specifically,
731 we want to bound:

$$732 \quad \left| \log(\bar{\pi}_{(m,T)}(a|s)) - \log(\pi_{\theta_{(m,T-1)}^i}(a|s)) \right| \leq \Delta_Q,$$

733 where Δ_Q involves the difference between the soft-Q values $\bar{Q}_{(m,T-1)}^{\text{soft}}(s, a)$ and
734 $Q_{r_{\theta_{(m,T-2)}^i}, \pi_{\theta_{(m,T-2)}^i}}^{\text{soft}}(s, a)$.

735 Recall that the policy is proportional to the exponential of the soft-Q value in equation 6. Thus, we
736 can write:

$$737 \quad \log(\bar{\pi}_{(m,T)}(a|s)) = \log \left(\frac{\exp(\bar{Q}_{m,T-1}^{\text{soft}}(s, a))}{\sum_{\tilde{a}} \exp(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a}))} \right)$$

$$738 \quad = \bar{Q}_{(m,T-1)}^{\text{soft}}(s, a) - \log \left(\sum_{\tilde{a}} \exp(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a})) \right) \quad (31)$$

739 Since $\log(\pi_{\theta_{(m,T-1)}^i}(a|s))$ is the optimal policy under reward parameter $\theta_{(m,T-1)}^i$, according to
740 (Haarnoja et al., 2017), it has the form

$$741 \quad \pi_{\theta_{(m,T-1)}^i}(a|s) = \frac{\exp(Q_{r_{\theta_{(m,T-2)}^i}, \pi_{\theta_{(m,T-2)}^i}}^{\text{soft}}(a|s)(s, a))}{\sum_{\tilde{a}} \exp(Q_{r_{\theta_{(m,T-2)}^i}, \pi_{\theta_{(m,T-2)}^i}}^{\text{soft}}(a|s)(s, \tilde{a}))} \quad (32)$$

Similarly, we have:

$$\log(\pi_{\theta^i(m,T-1)}(a|s)) = Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, a) - \log\left(\sum_{\tilde{a}} \exp\left(Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a})\right)\right) \quad (33)$$

Subtracting the two expressions and use the triangle inequality, we can bound the absolute value of the difference by the sum of the absolute values:

$$\begin{aligned} & \left| \log(\bar{\pi}_{(m,T)}(a|s)) - \log(\pi_{\theta^i(m,T-1)}(a|s)) \right| \\ &= \left| \left[\bar{Q}_{(m,T-1)}^{\text{soft}}(s, a) - \log\left(\sum_{\tilde{a}} \exp\left(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a})\right)\right) \right] \right. \\ & \quad \left. - \left[Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, a) - \log\left(\sum_{\tilde{a}} \exp\left(Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a})\right)\right) \right] \right| \\ &= \left| \bar{Q}_{(m,T-1)}^{\text{soft}}(s, a) - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, a) \right. \\ & \quad \left. - \left[\log\left(\sum_{\tilde{a}} \exp\left(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a})\right)\right) - \log\left(\sum_{\tilde{a}} \exp\left(Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a})\right)\right) \right] \right| \\ &\leq \left| \bar{Q}_{(m,T-1)}^{\text{soft}}(s, a) - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, a) \right| \\ & \quad + \left| \log\left(\sum_{\tilde{a}} \exp\left(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a})\right)\right) - \log\left(\sum_{\tilde{a}} \exp\left(Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a})\right)\right) \right| \quad (34) \end{aligned}$$

The second term in equation 34 involves the difference of logarithms of sums. We can bound it using properties of logarithms and the maximum difference of the soft-Q values.

We utilize the following inequality (as referenced in Equation 47 of (Zeng et al., 2022)):

$$\left| \log\left(\sum_{\tilde{a}} \exp(Q_1(s, \tilde{a}))\right) - \log\left(\sum_{\tilde{a}} \exp(Q_2(s, \tilde{a}))\right) \right| \leq \max_{\tilde{a}} |Q_1(s, \tilde{a}) - Q_2(s, \tilde{a})| \quad (35)$$

Applying equation 35, we get:

$$\begin{aligned} & \left| \log\left(\sum_{\tilde{a}} \exp\left(\bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a})\right)\right) - \log\left(\sum_{\tilde{a}} \exp\left(Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a})\right)\right) \right| \\ &\leq \max_{\tilde{a}} \left| \bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a}) - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a}) \right| \quad (36) \end{aligned}$$

Combining the results from equation 34 and equation 36:

$$\begin{aligned} & \left| \log(\bar{\pi}_{(m,T)}(a|s)) - \log(\pi_{\theta^i(m,T-1)}(a|s)) \right| \\ &\leq \left| \bar{Q}_{(m,T-1)}^{\text{soft}}(s, a) - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, a) \right| \\ & \quad + \max_{\tilde{a}} \left| \bar{Q}_{(m,T-1)}^{\text{soft}}(s, \tilde{a}) - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}(s, \tilde{a}) \right| \quad (37) \end{aligned}$$

810 Taking the infinity norm on equation 37 gives:

$$811 \quad \|\log(\bar{\pi}_{(m,T)}) - \log(\pi_{\theta^i_{(m,T-1)}})\|_\infty \leq 2\|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty \quad (38)$$

812 By the definition of aggregation in equation 10, we have

$$813 \quad \bar{Q}_{(m,T-1)}^{\text{soft}} = \frac{1}{N} \sum_{j=1}^N Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \quad (39)$$

814 Plug above definition into equation 38 and by triangle inequality:

$$815 \quad \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty = \left\| \sum_{j=1}^N \frac{1}{N} (Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}) \right\|_\infty$$

$$816 \quad \leq \frac{1}{N} \sum_{j=1}^N \|Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty \quad (40)$$

817 Therefore, we move to analyse $\|Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty$, which is the dif-

818 ference of soft-Q values between two different local nodes, one under policy estimation, and the
819 other under optimal policy. Looking back to the time right after last aggregation, where all local
820 servers have the same reward parameter $\bar{\theta}_{(m-1,T)}$, we could further bound this difference using the
821 difference of reward parameters, since the difference of local reward parameters are introduced by
822 the local increment at each internal iteration except for the aggregation round.

823 We start by decomposing this difference into three terms and use the triangle inequality, we bound
824 the sum :

$$825 \quad \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right\|_\infty$$

$$826 \quad = \left\| \left(Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \right) \right.$$

$$827 \quad + \left(Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^j_{(m,0)}}, \pi_{\theta^j_{(m,0)}}}^{\text{soft}} \right)$$

$$828 \quad + \left. \left(Q_{r_{\theta^i_{(m,0)}}, \pi_{\theta^i_{(m,0)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right) \right\|_\infty \quad (41)$$

$$829 \quad \leq \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \right\|_\infty$$

$$830 \quad + \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} - Q_{r_{\theta^j_{(m,0)}}, \pi_{\theta^j_{(m,0)}}}^{\text{soft}} \right\|_\infty$$

$$831 \quad + \left\| Q_{r_{\theta^i_{(m,0)}}, \pi_{\theta^i_{(m,0)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right\|_\infty$$

832 The first term is the difference between the soft-Q values under the same reward parameter $\theta^j_{(m,T-2)}$
833 but different policies $\pi_{\theta^j_{(m,T-2)}}$ and $\pi_{\theta^j_{(m,T-2)}}$. The second term is the difference due to the change
834 in reward parameters from $\theta^j_{(m,T-2)}$ to $\theta^j_{(m,0)}$, with corresponding optimal policies, and the third
835 term is similar to the second term but for node i , comparing $\theta^i_{(m,0)}$ and $\theta^i_{(m,T-2)}$. We are utilizing
836 the fact that $\theta^j_{(m,0)} = \theta^i_{(m,0)}$ since they are initialized after previous aggregation.

Applying equation 15 to equation 41, we have:

$$\begin{aligned}
& \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} \\
& \leq \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} \\
& \quad + L_q \left\| \theta^j_{(m,T-2)} - \theta^j_{(m,0)} \right\| + L_q \left\| \theta^i_{(m,0)} - \theta^i_{(m,T-2)} \right\|
\end{aligned} \tag{42}$$

Next, we express the differences in reward parameters in terms of gradient updates. According to equation 9:

$$\theta^j_{(m,T-2)} = \theta^j_{(m,0)} + \alpha \sum_{k=0}^{T-3} g^j_{(m,k)} \tag{43}$$

$$\theta^i_{(m,T-2)} = \theta^i_{(m,0)} + \alpha \sum_{k=0}^{T-3} g^i_{(m,k)} \tag{44}$$

Substituting back into our equation 42:

$$\begin{aligned}
& \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} \\
& \leq \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} + L_q \alpha \left\| \sum_{k=0}^{T-3} g^j_{(m,k)} \right\| + L_q \alpha \left\| \sum_{k=0}^{T-2} g^i_{(m,k)} \right\|
\end{aligned} \tag{45}$$

According to equation (56) in (Zeng et al., 2022)), the gradients are bounded as:

$$\|g^i_{(m,k)}\| \leq 2L_q \tag{46}$$

we can further bound the sums in equation 45:

$$\begin{aligned}
& \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} \\
& \leq \left\| Q_{r_{\theta^j_{(m,T-2)}}, \pi^j_{(m,T-2)}}^{\text{soft}} - Q_{r_{\theta^j_{(m,T-2)}}, \pi_{\theta^j_{(m,T-2)}}}^{\text{soft}} \right\|_{\infty} + 4(T-2)\alpha L_q^2
\end{aligned} \tag{47}$$

Now the difference is bounded by the difference between Q-values under reward parameter $\theta^j_{(m,T-2)}$ with respect to the previous approximated policy $\pi^j_{(m,T-2)}$ and the optimal policy $\pi_{\theta^j_{(m,T-2)}}$, plus some error terms. We come back from comparing Q-values across different nodes to evaluating the soft-Q value approximation within a single node.

By equation 57 in (Zeng et al., 2022):

$$\begin{aligned}
& \left\| Q_{r_{\theta^i_{(m,k)}}, \pi^i_{(m,k)}}^{\text{soft}} - Q_{r_{\theta^i_{(m,k)}}, \pi_{\theta^i_{(m,k)}}}^{\text{soft}} \right\|_{\infty} \leq \gamma \left\| Q_{r_{\theta^i_{(m,k-1)}}, \pi^i_{(m,k-1)}}^{\text{soft}} - Q_{r_{\theta^i_{(m,k-1)}}, \pi_{\theta^i_{(m,k-1)}}}^{\text{soft}} \right\|_{\infty} \\
& \quad + 4\alpha L_q^2, \quad 1 \leq k \leq T-1, \quad m \in \mathbb{N}, \quad \forall i
\end{aligned} \tag{48}$$

Above inequality provides a bounds of the local Q-value using the previous Q-value times a contraction factor γ plus some extra term. We could use it to compare the aggregated Q-value at m-th outer round with the aggregated Q-value at $m-1$ -th outer round:

$$\begin{aligned}
& \|Q_{r_{\theta^i(m,k)}, \pi^i(m,k)}^{\text{soft}} - Q_{r_{\theta^i(m,k)}, \pi_{\theta^i(m,k)}}^{\text{soft}}\|_{\infty} \\
& \leq \gamma^k \|Q_{r_{\theta^i(m,0)}, \pi^i(m,0)}^{\text{soft}} - Q_{r_{\theta^i(m,0)}, \pi_{\theta^i(m,0)}}^{\text{soft}}\|_{\infty} + \sum_{i=0}^{k-1} \gamma^i \cdot 4\alpha L_q^2 \\
& = \gamma^k \|Q_{r_{\theta^i(m,0)}, \pi^i(m,0)}^{\text{soft}} - Q_{r_{\theta^i(m,0)}, \pi_{\theta^i(m,0)}}^{\text{soft}}\|_{\infty} + \frac{1-\gamma^k}{1-\gamma} \cdot 4\alpha L_q^2, \quad m \in \mathbb{N}, \quad 1 \leq k \leq T-1
\end{aligned} \tag{49}$$

Apply equation 49 to equation 47:

$$\begin{aligned}
& \|Q_{r_{\theta^j(m,T-2)}, \pi^j(m,T-2)}^{\text{soft}} - Q_{r_{\theta^j(m,T-2)}, \pi_{\theta^j(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& \leq \|Q_{r_{\theta^j(m,T-2)}, \pi^j(m,T-2)}^{\text{soft}} - Q_{r_{\theta^j(m,T-2)}, \pi_{\theta^j(m,T-2)}}^{\text{soft}}\|_{\infty} + 4(T-2)\alpha L_q^2 \\
& \leq \gamma^{T-2} \|Q_{r_{\theta^j(m,0)}, \pi^j(m,0)}^{\text{soft}} - Q_{r_{\theta^j(m,0)}, \pi_{\theta^j(m,0)}}^{\text{soft}}\|_{\infty} + \frac{1-\gamma^{T-2}}{1-\gamma} \cdot 4\alpha L_q^2 + 4(T-2)\alpha L_q^2 \\
& = \gamma^{T-2} \|Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty} + 4\alpha \left(\frac{1-\gamma^{T-2}}{1-\gamma} + T-2 \right) L_q^2
\end{aligned} \tag{50}$$

Plug equation 50 into equation 40:

$$\begin{aligned}
& \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& \leq \frac{1}{N} \sum_{j=1}^N \|Q_{r_{\theta^j(m,T-2)}, \pi^j(m,T-2)}^{\text{soft}} - Q_{r_{\theta^i(m,T-1)}, \pi_{\theta^i(m,T-1)}}^{\text{soft}}\|_{\infty} \\
& \leq \frac{1}{N} \sum_{j=1}^N \left[\gamma^{T-2} \|Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty} + 4\alpha \left(2 \cdot \frac{1-\gamma^{T-2}}{1-\gamma} + T-2 \right) L_q^2 \right] \\
& = \gamma^{T-2} \|Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty} + 4\alpha \left(\frac{1-\gamma^{T-2}}{1-\gamma} + T-2 \right) L_q^2
\end{aligned} \tag{51}$$

We further analyze $\|Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty}$ and use triangle inequality to decompose it into three parts to acquire the same form of Q-value difference in the previous outer round:

$$\begin{aligned}
& \|Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty} \\
& = \|(Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - \bar{Q}_{(m-1,T-1)}^{\text{soft}}) + (\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m-1,T-1)}, \pi_{\theta^i(m-1,T-1)}}^{\text{soft}}) \\
& \quad + (Q_{r_{\theta^i(m-1,T-1)}, \pi_{\theta^i(m-1,T-1)}}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}})\|_{\infty} \\
& \leq \|(Q_{r_{\bar{\theta}(m-1,T)}, \bar{\pi}(m-1,T)}^{\text{soft}} - \bar{Q}_{(m-1,T-1)}^{\text{soft}})\|_{\infty} + \|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m-1,T-1)}, \pi_{\theta^i(m-1,T-1)}}^{\text{soft}}\|_{\infty} \\
& \quad + \|Q_{r_{\theta^i(m-1,T-1)}, \pi_{\theta^i(m-1,T-1)}}^{\text{soft}} - Q_{r_{\bar{\theta}(m-1,T)}, \pi_{\bar{\theta}(m-1,T)}}^{\text{soft}}\|_{\infty}
\end{aligned} \tag{52}$$

We first bound the first term in equation 52 by introducing an middle term and triangle inequality:

$$\begin{aligned}
& \|Q_{r_{\bar{\theta}}(m-1,T),\bar{\pi}(m-1,T)}^{\text{soft}} - \bar{Q}_{(m-1,T-1)}^{\text{soft}}\|_{\infty} \\
&= \|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - \frac{1}{N} \sum_{j=1}^N Q_{r_{\theta^j}(m-1,T-2),\pi^j(m-1,T-2)}^{\text{soft}}\|_{\infty} \\
&= \|(Q_{r_{\bar{\theta}}(m-1,T),\bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}}) \\
&\quad + (Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}} - \frac{1}{N} \sum_{j=1}^N Q_{r_{\theta^j}(m-1,T-2),\pi^j(m-1,T-2)}^{\text{soft}})\| \\
&\leq \|Q_{r_{\bar{\theta}}(m-1,T),\bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}}\|_{\infty} \\
&\quad + \frac{1}{N} \sum_{j=1}^N \|Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\theta^j}(m-1,T-2),\pi^j(m-1,T-2)}^{\text{soft}}\|_{\infty}
\end{aligned} \tag{53}$$

For the first term in equation 53, we leverage Lemma 7 in (Zeng et al., 2022), which states:

$$|Q_{r_{\theta_1},\pi}^{\text{soft}} - Q_{r_{\theta_2},\pi}^{\text{soft}}| \leq L_q \|\theta_1 - \theta_2\|, \forall \pi, \forall \theta_1, \theta_2, \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \tag{54}$$

Then we could further bound the difference between θ 's using equation 9 and equation 46:

$$\begin{aligned}
& \|Q_{r_{\bar{\theta}}(m-1,T),\bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}}\|_{\infty} \\
&\leq L_q \|\bar{\theta}_{(m-1,T)} - \theta_{(m-1,T-2)}^j\|_{\infty} \\
&= L_q \|\frac{1}{N} \sum_{j=1}^N \theta_{(m-1,T-1)}^j - \theta_{(m-1,T-2)}^j\| \\
&\leq \frac{L_q}{N} \sum_{j=1}^N \|\theta_{(m-1,T-1)}^j - \theta_{(m-1,T-2)}^j\| \\
&= \frac{L_q}{N} \sum_{j=1}^N \|(\theta_{(m-1,T-1)}^j - \theta_{(m-1,0)}^j) + (\theta_{(m-1,0)}^j - \theta_{(m-1,T-2)}^j)\| \\
&\leq \frac{L_q}{N} \sum_{j=1}^N (\|\theta_{(m-1,T-1)}^j - \theta_{(m-1,0)}^j\| + \|\theta_{(m-1,0)}^j - \theta_{(m-1,T-2)}^j\|) \\
&= \frac{L_q}{N} \sum_{j=1}^N \left(\alpha \left\| \sum_{k=0}^{T-2} g_{(m,k)}^j \right\| + \alpha \left\| \sum_{k=0}^{T-3} g_{(m,k)}^j \right\| \right) \\
&\leq 2(2T-3)\alpha L_q^2
\end{aligned}$$

For the second term in equation 53, by Lemma 2:

$$\begin{aligned}
& \|Q_{r_{\theta^j}(m-1,T-2),\bar{\pi}(m-1,T)}^{\text{soft}} - Q_{r_{\theta^j}(m-1,T-2),\pi^j(m-1,T-2)}^{\text{soft}}\|_{\infty} \\
&\leq \frac{1-\gamma}{\gamma} \|\log(\bar{\pi}_{(m-1,T)}) - \log(\pi_{(m-1,T-2)}^j)\|_{\infty}
\end{aligned} \tag{55}$$

We get result similar to equation 38 using similar techniques and decompose equation 55 using triangle inequality:

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036

$$\begin{aligned}
& \|\log(\bar{\pi}_{(m-1,T)}) - \log(\pi_{(m-1,T-2)}^j)\|_\infty \\
& \leq 2\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-3)}}, \pi_{\theta^i_{(m-1,T-3)}}}^{\text{soft}}\|_\infty \\
& = 2\|(\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}}) + (Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-3)}}, \pi_{\theta^i_{(m-1,T-3)}}}^{\text{soft}})\|_\infty \\
& \leq 2(\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}}\|_\infty + \|Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-3)}}, \pi_{\theta^i_{(m-1,T-3)}}}^{\text{soft}}\|_\infty)
\end{aligned} \tag{56}$$

1037
1038
1039
1040

The last term could be controlled according to what we did for the last two terms in equation 41:

1041
1042
1043
1044

$$\|Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-3)}}, \pi_{\theta^i_{(m-1,T-3)}}}^{\text{soft}}\|_\infty \leq 2\alpha L_q^2 \tag{57}$$

1045

Plugging above results into equation 52, we have:

1046
1047
1048
1049
1050
1051
1052
1053
1054

$$\begin{aligned}
& \|Q_{r_{\bar{\theta}_{(m-1,T)}}, \bar{\pi}_{(m-1,T)}}^{\text{soft}} - Q_{r_{\bar{\theta}_{(m-1,T)}}, \pi_{\bar{\theta}_{(m-1,T)}}}^{\text{soft}}\|_\infty \\
& \leq \frac{1-\gamma}{\gamma} \|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}}\|_\infty + 2(2T-3)\alpha L_q^2 + \frac{1-\gamma}{\gamma} \cdot 4\alpha L_q^2
\end{aligned} \tag{58}$$

1055

We further plug equation 58 to equation 51:

1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068

$$\begin{aligned}
& \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty \\
& \leq \gamma^{T-2} \|Q_{r_{\bar{\theta}_{(m-1,T)}}, \bar{\pi}_{(m-1,T)}}^{\text{soft}} - Q_{r_{\bar{\theta}_{(m-1,T)}}, \pi_{\bar{\theta}_{(m-1,T)}}}^{\text{soft}}\|_\infty + 4\alpha \left(\frac{1-\gamma^{T-2}}{1-\gamma} + T-2 \right) L_q^2 \\
& \leq (1-\gamma)\gamma^{T-1} \|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}}\|_\infty \\
& + \left(2\frac{1-\gamma^{T-2}}{1-\gamma} + (1-\gamma)^2\gamma^{T-2} + \frac{1-\gamma}{\gamma}(2T-3) + 2(T-2) \right) L_q^2
\end{aligned} \tag{59}$$

1069
1070

Summing the inequality from $m = 1$ to $m = M$ gives:

1071
1072
1073
1074
1075
1076
1077
1078
1079

$$\begin{aligned}
& \sum_{m=1}^M \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m,T-2)}}, \pi_{\theta^i_{(m,T-2)}}}^{\text{soft}}\|_\infty \\
& \leq (1-\gamma)\gamma^{T-1} \sum_{m=1}^M \|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}}^{\text{soft}}\|_\infty \\
& + M \cdot \left[2\frac{1-\gamma^{T-2}}{1-\gamma} + (1-\gamma)^2\gamma^{T-2} + \frac{1-\gamma}{\gamma}(2T-3) + 2(T-2) \right] L_q^2
\end{aligned} \tag{60}$$

Rearranging the inequality, it holds that:

$$\begin{aligned}
& (1 - (1 - \gamma)\gamma^{T-1}) \sum_{m=1}^M \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& \leq (1 - \gamma)\gamma^{T-1} \left(\|\bar{Q}_{(0,T-1)}^{\text{soft}} - Q_{r_{\theta^i(0,T-2)}, \pi_{\theta^i(0,T-2)}}^{\text{soft}}\|_{\infty} \right. \\
& \quad \left. - \|\bar{Q}_{(M,T-1)}^{\text{soft}} - Q_{r_{\theta^i(M,T-2)}, \pi_{\theta^i(M,T-2)}}^{\text{soft}}\|_{\infty} \right) \\
& \quad + M \cdot \left[2\frac{1 - \gamma^{T-2}}{1 - \gamma} + (1 - \gamma)^2\gamma^{T-2} + \frac{1 - \gamma}{\gamma}(2T - 3) + 2(T - 2) \right] L_q^2 \\
& \leq (1 - \gamma)\gamma^{T-1} \left(\|\bar{Q}_{(0,T-1)}^{\text{soft}} - Q_{r_{\theta^i(0,T-2)}, \pi_{\theta^i(0,T-2)}}^{\text{soft}}\|_{\infty} \right. \\
& \quad \left. + M \cdot \left[2\frac{1 - \gamma^{T-2}}{1 - \gamma} + (1 - \gamma)^2\gamma^{T-2} + \frac{1 - \gamma}{\gamma}(2T - 3) + 2(T - 2) \right] L_q^2 \right)
\end{aligned} \tag{61}$$

Dividing by $1 - (1 - \gamma)\gamma^{T-1}$ on both sides, we get

$$\begin{aligned}
& \sum_{m=1}^M \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& \leq \frac{(1 - \gamma)\gamma^{T-1}}{1 - (1 - \gamma)\gamma^{T-1}} \|\bar{Q}_{(0,T-1)}^{\text{soft}} - Q_{r_{\theta^i(0,T-2)}, \pi_{\theta^i(0,T-2)}}^{\text{soft}}\|_{\infty} \\
& \quad + M \cdot \frac{2\frac{1 - \gamma^{T-2}}{1 - \gamma} + (1 - \gamma)^2\gamma^{T-2} + \frac{1 - \gamma}{\gamma}(2T - 3) + 2(T - 2)}{1 - (1 - \gamma)\gamma^{T-1}} L_q^2
\end{aligned}$$

Denote $C_0 = \|\bar{Q}_{(0,T-1)}^{\text{soft}} - Q_{r_{\theta^i(0,T-2)}, \pi_{\theta^i(0,T-2)}}^{\text{soft}}\|_{\infty}$.

Dividing by M on both sides, we get

$$\begin{aligned}
& \frac{1}{M} \sum_{m=1}^M \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& \leq \frac{(1 - \gamma)\gamma^{T-1}}{(1 - (1 - \gamma)\gamma^{T-1})M} C_0 + \frac{2\frac{1 - \gamma^{T-2}}{1 - \gamma} + (1 - \gamma)^2\gamma^{T-2} + \frac{1 - \gamma}{\gamma}(2T - 3) + 2(T - 2)}{1 - (1 - \gamma)\gamma^{T-1}} L_q^2
\end{aligned} \tag{62}$$

Recall the step size is defined as $\alpha_{(m,t)} = \frac{\alpha_0}{(mT+t)^\sigma}$ where $\sigma > 0$. Then we have the following result:

$$\begin{aligned}
& \frac{1}{M} \sum_{m=1}^M \|\bar{Q}_{(m,T-1)}^{\text{soft}} - Q_{r_{\theta^i(m,T-2)}, \pi_{\theta^i(m,T-2)}}^{\text{soft}}\|_{\infty} \\
& = \mathcal{O}(M^{-1}\gamma^T) + \mathcal{O}(M^{-\sigma}T^{1-\sigma})
\end{aligned} \tag{63}$$

Going back to the convergence of policy approximation:

$$\begin{aligned}
& \frac{1}{M} \sum_{m=1}^M |\log(\bar{\pi}_{(m,T)}(a|s)) - \log(\pi_{\theta^i(m,T-1)}(a|s))|_{\infty} \\
& = \mathcal{O}(M^{-1}\gamma^T) + \mathcal{O}(M^{-\sigma}T^{1-\sigma})
\end{aligned} \tag{64}$$

B.2 CONVERGENCE OF THE GLOBAL REWARD PARAMETER $\bar{\theta}_{(m,T)}$

By the Lipschitz smooth property of the likelihood target equation 16, the definition of reward aggregation equation 12, and the reward parameter update rule equation 9:

$$\begin{aligned}
& L(\bar{\theta}_{(m,T)}) \\
& \geq L(\bar{\theta}_{(m-1,T)}) + \langle \nabla L(\bar{\theta}_{(m,T)}), \bar{\theta}_{(m,T)} - \bar{\theta}_{(m-1,T)} \rangle - \frac{L_c}{2} \|\bar{\theta}_{(m,T)} - \bar{\theta}_{(m-1,T)}\|^2 \\
& = L(\bar{\theta}_{(m-1,T)}) + \left\langle \nabla L(\bar{\theta}_{(m,T)}), \frac{1}{N} \sum_{j=1}^N (\theta_{(m,T-1)}^j - \theta_{(m,0)}^j) \right\rangle - \frac{L_c}{2} \left\| \frac{1}{N} \sum_{j=1}^N (\theta_{(m,T-1)}^j - \theta_{(m,0)}^j) \right\|^2 \\
& = L(\bar{\theta}_{(m-1,T)}) + \alpha \left\langle \nabla L(\bar{\theta}_{(m,T)}), \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j \right\rangle - \frac{L_c \alpha^2}{2} \left\| \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j \right\|^2
\end{aligned} \tag{65}$$

We compare $g_{(m,k)}^j$ with the true gradient of $L(\theta_{(m,k)}^j)$ and leverage the fact that $\|\nabla L(\theta)\|_\infty \leq 2L_q$:

$$\begin{aligned}
& L(\bar{\theta}_{(m,T)}) \\
& \geq L(\bar{\theta}_{(m-1,T)}) + \alpha \left\langle \nabla L(\bar{\theta}_{(m,T)}), \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j \right\rangle - \frac{L_c \alpha^2}{2} \left\| \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j \right\|^2 \\
& \geq L(\bar{\theta}_{(m-1,T)}) + \alpha \left\langle \nabla L(\bar{\theta}_{(m,T)}), \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j - (T-1) \nabla L(\bar{\theta}_{(m-1,T)}) \right\rangle \\
& \quad + \alpha(T-1) \|\nabla L(\bar{\theta}_{(m-1,T)})\|^2 - \frac{L_c \alpha^2}{2} \left\| \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} g_{(m,k)}^j \right\|^2 \\
& \geq L(\bar{\theta}_{(m-1,T)}) - 2\alpha L_q \cdot \left\| \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} [g_{(m,k)}^j - \nabla L(\bar{\theta}_{(m-1,T)})] \right\| \\
& \quad + \alpha(T-1) \|\nabla L(\bar{\theta}_{(m-1,T)})\|^2 - (T-1)^2 \cdot \frac{L_c L_q^2 \alpha^2}{2}
\end{aligned} \tag{66}$$

For $g^j(m, k)$ in equation 66, we evaluate its distance to $\nabla L(\theta_{(m,k)}^j)$ and also consider the distance between $\nabla L(\theta_{(m,k)}^j)$ and $\nabla L(\bar{\theta}_{(m,T)})$ with the help of triangle inequality:

$$\begin{aligned}
& L(\bar{\theta}_{(m,T)}) \\
& \geq L(\bar{\theta}_{(m-1,T)}) - 2\alpha L_q \cdot \left\| \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} [g_{(m,k)}^j - \nabla L(\theta_{(m,k)}^j) + \nabla L(\theta_{(m,k)}^j) - \nabla L(\theta_{(m,0)}^j)] \right\| \\
& \quad + \alpha(T-1) \|\nabla L(\bar{\theta}_{(m-1,T)})\|^2 - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \\
& \geq L(\bar{\theta}_{(m-1,T)}) - 2\alpha L_q \cdot \left[\frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} (\|g_{(m,k)}^j - \nabla L(\theta_{(m,k)}^j)\| + \|\nabla L(\theta_{(m,k)}^j) - \nabla L(\theta_{(m,0)}^j)\|) \right] \\
& \quad + \alpha(T-1) \|\nabla L(\bar{\theta}_{(m-1,T)})\|^2 - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2}
\end{aligned} \tag{67}$$

Taking expectation over both sides:

$$\begin{aligned}
& \mathbb{E}[L(\bar{\theta}_{(m,T)})] \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - 2\alpha L_q \cdot \left[\frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} (\mathbb{E}\|g_{(m,k)}^j - \nabla L(\theta_{(m,k)}^j)\| + \|\nabla L(\theta_{(m,k)}^j) - \nabla L(\theta_{(m,0)})\|) \right] \\
& \quad + \alpha(T-1) \|\nabla L(\bar{\theta}_{(m-1,T)})\|^2 - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2}
\end{aligned} \tag{68}$$

According to equation (62) and (63) in (Zeng et al., 2022), we have:

$$\mathbb{E}\|g_{(m,k)}^j - \nabla L(\theta_{(m,k)}^j)\| \leq 2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j} - Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j}\|]_{\infty} \tag{69}$$

Then, using the Lipschitz property of L in equation 16:

$$\begin{aligned}
& \mathbb{E}[L(\bar{\theta}_{(m,T)})] \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - 2\alpha L_q \cdot \frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} (2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j} - Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j}\|]_{\infty} \\
& \quad + \mathbb{E}[\|\nabla L(\theta_{(m,k)}^j) - \nabla L(\theta_{(m,0)}^j)\|]) + \alpha(T-1) \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - 2\alpha L_q \cdot \left[\frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} (2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j} - Q_{r_{\theta_{(m,k)}^j}^{\text{soft}}, \pi_{\theta_{(m,k)}^j}^j}\|]_{\infty} \right. \\
& \quad \left. + \mathbb{E}[L_c \|\theta_{(m,k)}^j - \theta_{(m,0)}^j\|]) + \alpha(T-1) \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \right]
\end{aligned} \tag{70}$$

Similar to equation 47, we have $\|\theta_{(m,k)}^j - \theta_{(m,0)}^j\| \leq 2k\alpha L_q$, applying equation 49 to equation 70, we have:

$$\begin{aligned}
& \mathbb{E}[L(\bar{\theta}_{(m,T)})] \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - 2\alpha L_q \cdot \left[\frac{1}{N} \sum_{j=1}^N \sum_{k=0}^{T-2} (2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \gamma^k \mathbb{E}[\|Q_{r_{\theta_{(m,0)}^j}^{\text{soft}}, \pi_{\theta_{(m,0)}^j}^j} - Q_{r_{\theta_{(m,0)}^j}^{\text{soft}}, \pi_{\theta_{(m,0)}^j}^j}\|]_{\infty}) \right. \\
& \quad \left. + 2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \frac{1-\gamma^k}{1-\gamma} \cdot 4\alpha L_q^2 + 2k\alpha L_c L_q + \alpha(T-1) \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \right] \\
& = \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - 2\alpha L_q \cdot \sum_{k=0}^{T-2} \gamma^k (2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\bar{\theta}_{(m-1,T)}}^{\text{soft}}, \bar{\pi}_{(m-1,T)}} - Q_{r_{\bar{\theta}_{(m-1,T)}}^{\text{soft}}, \bar{\pi}_{(m-1,T)}}\|]_{\infty}) \\
& \quad + 2L_q C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \frac{1-\gamma^k}{1-\gamma} \cdot 4\alpha L_q^2 + k\alpha L_c L_q + \alpha(T-1) \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \\
& = \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - \frac{4(1-\gamma^{T-1})}{1-\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\bar{\theta}_{(m-1,T)}}^{\text{soft}}, \bar{\pi}_{(m-1,T)}} - Q_{r_{\bar{\theta}_{(m-1,T)}}^{\text{soft}}, \bar{\pi}_{(m-1,T)}}\|]_{\infty} \\
& \quad - 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1-\frac{1-\gamma^{T-1}}{1-\gamma}}{1-\gamma} - T(T-1)\alpha^2 L_c L_q^2 + \alpha(T-1) \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] \\
& \quad - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2}
\end{aligned} \tag{71}$$

Plugging the result in equation 58, we have:

$$\begin{aligned}
& \mathbb{E}[L(\bar{\theta}_{(m,T)})] \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - \frac{4(1-\gamma^{T-1})}{1-\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|Q_{r_{\bar{\theta}_{(m-1,T)}}, \bar{\pi}_{(m-1,T)}}^{\text{soft}} - Q_{r_{\bar{\theta}_{(m-1,T)}}, \pi_{\bar{\theta}_{(m-1,T)}}^{\text{soft}}}\|_{\infty}] \\
& \quad - 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1-\frac{1-\gamma^{T-1}}{1-\gamma}}{1-\gamma} - T(T-1)\alpha^2 L_c L_q^2 + \alpha(T-1)\mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] \\
& \quad - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} \\
& \geq \mathbb{E}[L(\bar{\theta}_{(m-1,T)})] - \frac{4(1-\gamma^{T-1})}{\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \mathbb{E}[\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}^{\text{soft}}}\|_{\infty}] \\
& \quad + \alpha(T-1)\mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] - 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1-\frac{1-\gamma^{T-1}}{1-\gamma}}{1-\gamma} - T(T-1)\alpha^2 L_c L_q^2 \\
& \quad - \frac{(T-1)^2 L_c L_q^2 \alpha^2}{2} - \frac{4(1-\gamma^{T-1})}{1-\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot [2(2T-3)\alpha L_q^2 + \frac{1-\gamma}{\gamma} \cdot 4\alpha L_q^2]
\end{aligned} \tag{72}$$

Rearranging the inequality above and denote $C_1 = \frac{4(1-\gamma^{T-1})}{\gamma} L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}$, we obtain:

$$\begin{aligned}
& \alpha(T-1)\mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] \\
& \leq \alpha C_1 \mathbb{E}[\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}^{\text{soft}}}\|_{\infty}] + \mathbb{E}[L(\bar{\theta}_{(m,T)}) - L(\bar{\theta}_{(m-1,T)})] \\
& \quad + 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1-\frac{1-\gamma^{T-1}}{1-\gamma}}{1-\gamma} + \frac{(T-1)(3T-1)\alpha^2 L_c L_q^2}{2} \\
& \quad + \frac{4(1-\gamma^{T-1})}{1-\gamma} \alpha L_q^2 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot [2(2T-3)\alpha L_q^2 + \frac{1-\gamma}{\gamma} \cdot 4\alpha L_q^2] \\
& \leq \alpha C_1 \mathbb{E}[\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}^{\text{soft}}}\|_{\infty}] + \mathbb{E}[L(\bar{\theta}_{(m,T)}) - L(\bar{\theta}_{(m-1,T)})] \\
& \quad + 8\alpha L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} \cdot \frac{T-1}{1-\gamma} + \frac{(T-1)(3T-1)\alpha^2 L_c L_q^2}{2} \\
& \quad + \frac{16(T-1+\frac{1-\gamma}{\gamma})}{1-\gamma} \alpha^2 L_q^4 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}
\end{aligned} \tag{73}$$

Summing the inequality above from $m = 1$ to M and dividing both sides by $\alpha(T-1)M$, it holds that

$$\begin{aligned}
& \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1,T)})\|^2] \\
& \leq \frac{1}{M(T-1)} \sum_{m=1}^M \mathbb{E}[\|\bar{Q}_{(m-1,T-1)}^{\text{soft}} - Q_{r_{\theta^i_{(m-1,T-2)}}, \pi_{\theta^i_{(m-1,T-2)}}^{\text{soft}}}\|_{\infty}] + \mathbb{E}[\frac{L(\bar{\theta}_{(m,T)}) - L(\bar{\theta}_{(0,T)})}{\alpha(T-1)M}] \\
& \quad + \frac{1}{M} \frac{8}{1-\gamma} L_q^3 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|} + \frac{(3T-1)\alpha L_c L_q^2}{M} \frac{1}{2} + \frac{16\alpha(1+\frac{1-\gamma}{\gamma})}{M} \frac{L_q^4 C_d \sqrt{|\mathcal{S}| \cdot |\mathcal{A}|}}{1-\gamma}
\end{aligned} \tag{74}$$

Since $L(\bar{\theta}_{(m,T)})$ is negative and $L(\bar{\theta}_{(0,T)})$ is bounded constant, we plug equation 63 into equation 74 and get:

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E}[\|\nabla L(\bar{\theta}_{(m-1, T)})\|^2] = \mathcal{O}(M^{-1}) + \mathcal{O}(M^{-\sigma} T^{-\sigma}) + \mathcal{O}(M^{-1-\sigma} T^{1-\sigma}) \quad (75)$$

C EVALUATION

We present the details of the experiment setup and show more convergence plots in MuJoCo tasks.

C.1 EXPERIMENT SETUP

For f-IRL, we utilize the official implementation available at <https://github.com/twni2016/f-IRL>, which also includes implementations for BC and GAIL. The official implementation of ML-IRL can be found at <https://github.com/Cloud0723/ML-IRL>.

To ensure a fair comparison, we use SAC as the base RL algorithm for our F-ML-IRL approach as well as for all baselines, and Adam as the optimizer. Both the Q-network and policy network are configured as 64×64 MLPs with ReLU activation functions, and the learning rate is set to 1×10^{-3} .

For the Ant and Humanoid environments, the reward function is parameterized by a 128×128 MLP with ReLU activation, while for HalfCheetah, Hopper, and Walker2d, a 64×64 MLP with ReLU activation is used. The learning rate for the reward parameter is 1×10^{-4} for Hopper and 1×10^{-3} for the other environments.

At each iteration, we sample 10 trajectories from the current local policy estimate and compare them with the expert demonstration to update the reward parameter.

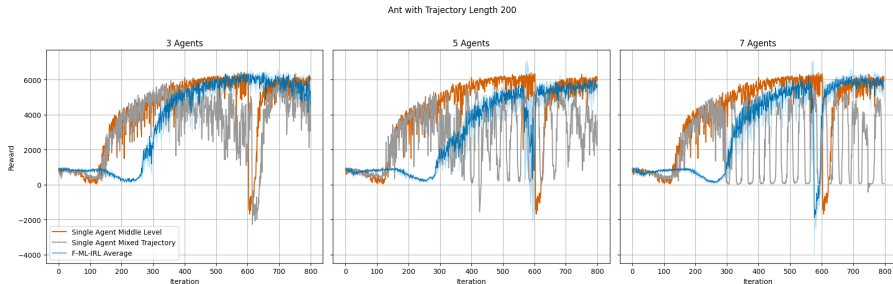
The reward levels of the expert demonstrations are shown in Table C.1. For 3 agents, we use Data 3, 4, and 5; for 5 agents, Data 2, 3, 4, 5, and 6 are used. For 7 agents, all 7 data sets are distributed across different local clients.

Environment	Data 1	Data 2	Data 3	Data 4	Data 5	Data 6	Data 7
Ant	5465.10	5544.83	5699.96	5758.39	5820.41	5927.86	6035.14
HalfCheetah	12831.84	12973.62	13045.36	13187.47	13236.31	13328.38	13434.40
Hopper	3122.05	3217.36	3305.78	3424.81	3553.03	3603.60	3709.89
Humanoid	4934.42	5074.53	5134.65	5297.35	5345.38	5420.32	5501.73
Walker2d	4801.82	4976.41	5081.29	5193.50	5220.25	5379.48	5440.2

Table 2: The reward levels of 7 expert demonstration datasets that are used in our experiments across 5 MuJoCo tasks. We distribute these non-iid datasets to the clients in our experiments.

C.2 CONVERGENCE PLOTS

We provide supplementary plots in other settings (different environment and trajectory length) here to show the convergence of F-ML-IRL compared with ML-IRL in two centralized learning data cases:



1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

