
Generative Adversarial Imitation Learning with Neural Network Parameterization: Global Optimality and Convergence Rate

Yufeng Zhang¹ Qi Cai¹ Zhuoran Yang² Zhaoran Wang¹

Abstract

Generative adversarial imitation learning (GAIL) demonstrates tremendous success in practice, especially when combined with neural networks. Different from reinforcement learning, GAIL learns both policy and reward function from expert (human) demonstration. Despite its empirical success, it remains unclear whether GAIL with neural networks converges to the globally optimal solution. The major difficulty comes from the nonconvex-nonconcave minimax optimization structure. To bridge the gap between practice and theory, we analyze a gradient-based algorithm with alternating updates and establish its sublinear convergence to the globally optimal solution. To the best of our knowledge, our analysis establishes the global optimality and convergence rate of GAIL with neural networks for the first time.

1. Introduction

The goal of imitation learning (IL) is to learn to perform a task based on expert demonstration (Ho & Ermon, 2016). In contrast to reinforcement learning (RL), the agent only has access to the expert trajectories but not the rewards. The most straightforward approach of IL is behavioral cloning (BC) (Pomerleau, 1991). BC treats IL as the supervised learning problem of predicting the actions based on the states. Despite its simplicity, BC suffers from the compounding errors caused by covariate shift (Ross et al., 2011; Ross & Bagnell, 2010). Another approach of IL is inverse reinforcement learning (IRL) (Russell, 1998; Ng & Russell, 2000; Levine & Koltun, 2012; Finn et al., 2016), which jointly learns the reward function and the corresponding

optimal policy. IRL formulates IL as a bilevel optimization problem. Specifically, IRL solves an RL subproblem given a reward function at the inner level and searches for the reward function which makes the expert policy optimal at the outer level. However, IRL is computationally inefficient as it requires fully solving an RL subproblem at each iteration of the outer level. Moreover, the desired reward function may be nonunique. To address such issues of IRL, (Ho & Ermon, 2016) propose generative adversarial imitation learning (GAIL), which searches for the optimal policy without fully solving an RL subproblem given a reward function at each iteration. GAIL solves IL via minimax optimization with alternating updates. In particular, GAIL alternates between (i) minimizing the discrepancy in expected cumulative reward between the expert policy and the learned policy and (ii) maximizing such a discrepancy over the reward function class. Such an alternating update scheme mirrors the training of generative adversarial networks (GANs) (Goodfellow et al., 2014; Arjovsky et al., 2017), where the learned policy acts as the generator while the reward function acts as the discriminator.

Incorporated with neural networks, which parameterize the learned policy and the reward function, GAIL achieves significant empirical success in challenging applications, such as natural language processing (Yu et al., 2016), autonomous driving (Kuefler et al., 2017), human behavior modeling (Merel et al., 2017), and robotics (Tai et al., 2018). Despite its empirical success, GAIL with neural networks remains less understood in theory. The major difficulty arises from the following aspects: (i) GAIL involves minimax optimization, while the existing analysis of policy optimization with neural networks (Anthony & Bartlett, 2009; Liu et al., 2019; Bhandari & Russo, 2019; Wang et al., 2019) only focuses on a minimization or maximization problem. (ii) GAIL with neural networks is nonconvex-nonconcave, and therefore, the existing analysis of convex-concave optimization with alternating updates is not applicable (Nesterov, 2013). There is an emerging body of literature (Rafique et al., 2018; Zhang et al., 2019b) that casts nonconvex-nonconcave optimization as bilevel optimization, where the inner level is solved to a high precision as in IRL. However, such analysis is not applicable to GAIL as it involves alternating updates.

¹Department of Industrial Engineering and Management Sciences, Northwestern University, Evanston, IL 60208, USA
²Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544, USA. Correspondence to: Zhaoran Wang <zhaoranwang@gmail.com>.

In this paper, we bridge the gap between practice and theory by establishing the global optimality and convergence of GAIL with neural networks. Specifically, we parameterize the learned policy and the reward function with two-layer neural networks and consider solving GAIL by alternatively updating the learned policy via a step of natural policy gradient (Kakade, 2002; Peters & Schaal, 2008) and the reward function via a step of gradient ascent. In particular, we parameterize the state-action value function (also known as the Q-function) with a two-layer neural network and apply a variant of the temporal difference algorithm (Sutton & Barto, 2018) to solve the policy evaluation subproblem in natural policy gradient. We prove that the learned policy $\bar{\pi}$ converges to the expert policy π_E at a $1/\sqrt{T}$ rate in the \mathcal{R} -distance (Chen et al., 2020), which is defined as $\mathbb{D}_{\mathcal{R}}(\pi_E, \bar{\pi}) = \max_{r \in \mathcal{R}} J(\pi_E; r) - J(\bar{\pi}; r)$. Here $J(\pi; r)$ is the expected cumulative reward of a policy π given a reward function $r(s, a)$ and \mathcal{R} is the reward function class. The core of our analysis is constructing a potential function that tracks the \mathcal{R} -distance. Such a rate of convergence implies that the learned policy $\bar{\pi}$ (approximately) outperforms the expert policy π_E given any reward function $r \in \mathcal{R}$ within a finite number of iterations T . In other words, the learned policy $\bar{\pi}$ is globally optimal. To the best of our knowledge, our analysis establishes the global optimality and convergence of GAIL with neural networks for the first time. It is worth mentioning that our analysis is straightforwardly applicable to linear and tabular settings, which, however, are not our focus.

Related works. Our work extends an emerging body of literature on RL with neural networks (Xu et al., 2019a; Zhang et al., 2019a; Bhandari & Russo, 2019; Liu et al., 2019; Wang et al., 2019; Agarwal et al., 2019) to IL. This line of research analyzes the global optimality and convergence of policy gradient for solving RL, which is a minimization or maximization problem. In contrast, we analyze GAIL, which is a minimax optimization problem.

Our work is also related to the analysis of apprenticeship learning (Syed et al., 2008) and GAIL (Cai et al., 2019a; Chen et al., 2020). (Syed et al., 2008) analyze the convergence and generalization of apprenticeship learning. They assume the state space to be finite, and thus, do not require function approximation for the policy and the reward function. In contrast, we assume the state space to be infinite and employ function approximation based on neural networks. (Cai et al., 2019a) study the global optimality and convergence of GAIL in the setting of linear-quadratic regulators. In contrast, our analysis handles general MDPs without restrictive assumptions on the transition kernel and the reward function. (Chen et al., 2020) study the convergence and generalization of GAIL in the setting of general MDPs. However, they only establish the convergence to a stationary point. In contrast, we establish the global optimality of

GAIL.

Notations. Let $[n] = \{1, \dots, n\}$ for $n \in \mathbb{N}_+$ and $[m : n] = \{m, m + 1, \dots, n\}$ for $m < n$. Also, let $N(\mu, \Sigma)$ be the Gaussian distribution with mean μ and covariance Σ . We denote by $\mathcal{P}(\mathcal{X})$ the set of all probability measures over the space \mathcal{X} . For a function $f : \mathcal{X} \rightarrow \mathbb{R}$, a constant $p \geq 1$, and a probability measure $\mu \in \mathcal{P}(\mathcal{X})$, we denote by $\|f\|_{p, \mu} = (\int_{\mathcal{X}} |f(x)|^p d\mu(x))^{1/p}$ the $L_p(\mu)$ norm of the function f . For any two functions $f, g : \mathcal{X} \rightarrow \mathbb{R}$, we denote by $\langle f, g \rangle_{\mathcal{X}} = \int_{\mathcal{X}} f(x) \cdot g(x) dx$ the inner product on the space \mathcal{X} .

2. Background

In this section, we introduce reinforcement learning (RL) and generative adversarial imitation learning (GAIL).

2.1. Reinforcement Learning

We consider a Markov decision process (MDP) $(\mathcal{S}, \mathcal{A}, r, P, \rho, \gamma)$. Here $\mathcal{S} \subseteq \mathbb{R}^{d_1}$ is the state space, $\mathcal{A} \subseteq \mathbb{R}^{d_2}$ is the action space, which is assumed to be finite throughout this paper, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ is the transition kernel, $\rho \in \mathcal{P}(\mathcal{S})$ is the initial state distribution, and $\gamma \in (0, 1)$ is the discount factor. Without loss of generality, we assume that $\mathcal{S} \times \mathcal{A}$ is compact and that $\|(s, a)\|_2 \leq 1$ for any $(s, a) \in \mathcal{S} \times \mathcal{A} \subseteq \mathbb{R}^d$, where $d = d_1 + d_2$. An agent following a policy $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ interacts with the environment in the following manner. At the state $s_t \in \mathcal{S}$, the agent takes the action $a_t \in \mathcal{A}$ with probability $\pi(a_t | s_t)$ and receives the reward $r_t = r(s_t, a_t)$. The environment then transits into the next state s_{t+1} with probability $P(s_{t+1} | s_t, a_t)$. Given a policy π and a reward function $r(s, a)$, we define the state-action value function $Q_r^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as follows,

$$\begin{aligned} Q_r^\pi(s, a) & \quad (2.1) \\ & = \mathbb{E}_\pi \left[(1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot r(s_t, a_t) \mid s_0 = s, a_0 = a \right]. \end{aligned}$$

Here the expectation \mathbb{E}_π is taken with respect to $a_t \sim \pi(\cdot | s_t)$ and $s_{t+1} \sim P(\cdot | s_t, a_t)$. Correspondingly, we define the state value function $V_r^\pi : \mathcal{S} \rightarrow \mathbb{R}$ and the advantage function $A_r^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as follows,

$$\begin{aligned} V_r^\pi(s) & = \mathbb{E}_{a \sim \pi(\cdot | s)} [Q_r^\pi(s, a)], \\ A_r^\pi(s, a) & = Q_r^\pi(s, a) - V_r^\pi(s). \end{aligned}$$

The goal of RL is to maximize the following expected cumulative reward,

$$J(\pi; r) = \mathbb{E}_{s \sim \rho} [V_r^\pi(s)]. \quad (2.2)$$

The policy π induces a state visitation measure $d_\pi \in \mathcal{P}(\mathcal{S})$ and a state-action visitation measure $\nu_\pi \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$, which

take the forms of

$$\begin{aligned} d_\pi(s) &= (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}(s_t = s \mid s_0 \sim \rho, \pi), \\ \nu_\pi(s, a) &= d_\pi(s) \cdot \pi(a \mid s). \end{aligned} \quad (2.3)$$

It then holds that $J(\pi; r) = \mathbb{E}_{(s,a) \sim \nu_\pi} [r(s, a)]$. Meanwhile, we assume that the policy π induces a state stationary distribution ϱ_π over \mathcal{S} , which satisfies that

$$\varrho_\pi(s) = \mathbb{P}(s_{t+1} = s \mid s_t \sim \rho_\pi, a_t \sim \pi(\cdot \mid s_t)).$$

We denote by $\rho_\pi(s, a) = \varrho(s) \cdot \pi(a \mid s)$ the state-action stationary distribution over $\mathcal{S} \times \mathcal{A}$.

2.2. Generative Adversarial Imitation Learning

The goal of imitation learning (IL) is to learn a policy that outperforms the expert policy π_E based on the trajectory $\{(s_i^E, a_i^E)\}_{i \in [T_E]}$ of π_E . We denote by $\nu_E = \nu_{\pi_E}$ and $d_E = d_{\pi_E}$ the state-action and state visitation measures induced by the expert policy, respectively, and assume that the expert trajectory $\{(s_i, a_i)\}_{i \in [T_E]}$ is drawn from ν_E . To this end, we parameterize the policy and the reward function by π_θ for $\theta \in \mathcal{X}_\Pi$ and $r_\beta(s, a)$ for $\beta \in \mathcal{X}_R$, respectively, and solve the following minimax optimization problem known as GAIL (Ho & Ermon, 2016),

$$\begin{aligned} \min_{\theta \in \mathcal{X}_\Pi} \max_{\beta \in \mathcal{X}_R} L(\theta, \beta), \\ \text{where } L(\theta, \beta) = J(\pi_\theta; r_\beta) - J(\pi_E; r_\beta) - \lambda \cdot \psi(\beta). \end{aligned} \quad (2.4)$$

Here $J(\pi; r)$ is the expected cumulative reward defined in (2.2), $\psi : \mathcal{X}_R \rightarrow \mathbb{R}_+$ is the regularizer, and $\lambda \geq 0$ is the regularization parameter. Given a reward function class \mathcal{R} , we define the \mathcal{R} -distance between two policies π_1 and π_2 as follows,

$$\begin{aligned} \mathbb{D}_{\mathcal{R}}(\pi_1, \pi_2) &= \max_{r \in \mathcal{R}} J(\pi_1; r) - J(\pi_2; r) \\ &= \max_{r \in \mathcal{R}} \mathbb{E}_{\nu_{\pi_1}} [r(s, a)] - \mathbb{E}_{\nu_{\pi_2}} [r(s, a)]. \end{aligned} \quad (2.5)$$

When \mathcal{R} is the class of 1-Lipschitz functions, $\mathbb{D}_{\mathcal{R}}(\pi_1, \pi_2)$ is the Wasserstein-1 metric between the state-action visitation measures induced by π_1 and π_2 . However, $\mathbb{D}_{\mathcal{R}}(\pi_1, \pi_2)$ is not a metric in general. When $\mathbb{D}_{\mathcal{R}}(\pi_1, \pi_2) \leq 0$, the policy π_2 outperforms the policy π_1 for any reward function $r \in \mathcal{R}$. Such a notion of \mathcal{R} -distance is used in (Chen et al., 2020). We denote by $\mathcal{R}_\beta = \{r_\beta(s, a) \mid \beta \in \mathcal{X}_R\}$ the reward function class parameterized with β . Hence, the optimization problem in (2.4) minimizes the \mathcal{R}_β -distance $\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \pi_\theta)$ (up to the regularizer $\lambda \cdot \psi(\beta)$), which searches for a policy $\bar{\pi}$ that (approximately) outperforms the expert policy given any reward function $r_\beta \in \mathcal{R}_\beta$.

3. Algorithm

In this section, we introduce an algorithm with alternating updates for GAIL with neural networks, which employs natural policy gradient (NPG) to update the policy π_θ and gradient ascent to update the reward function $r_\beta(s, a)$.

3.1. Parameterization with Neural Networks

We define a two-layer neural network with rectified linear units (ReLU) as follows,

$$\begin{aligned} u_{W,b}(s, a) &= \frac{1}{\sqrt{m}} \sum_{l=1}^m b_l \cdot \mathbb{1} \{ (s, a)^\top [W]_l > 0 \} \cdot (s, a)^\top [W]_l \\ &= \sum_{l=1}^m [\phi_{W,b}(s, a)]_l^\top [W]_l. \end{aligned} \quad (3.1)$$

Here $m \in \mathbb{N}_+$ is the width of the neural network, $b = (b_1, \dots, b_m)^\top \in \mathbb{R}^m$ and $W = ([W]_1^\top, \dots, [W]_m^\top)^\top \in \mathbb{R}^{md}$ are the parameters, and $\phi_{W,b}(s, a) = ([\phi_{W,b}(s, a)]_1^\top, \dots, [\phi_{W,b}(s, a)]_m^\top)^\top \in \mathbb{R}^{md}$ is called the feature vector in the sequel, where

$$[\phi_{W,b}(s, a)]_l = \frac{b_l}{\sqrt{m}} \cdot \mathbb{1} \{ (s, a)^\top [W]_l > 0 \} (s, a). \quad (3.2)$$

It then holds that $u_{W,b}(s, a) = W^\top \phi_{W,b}(s, a)$. Note that the feature vector $\phi_{W,b}(s, a)$ depends on the parameters W and b . We consider the following random initialization,

$$b_l \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{-1, 1\}), \quad [W_0]_l \stackrel{\text{i.i.d.}}{\sim} N(0, I_d/d), \quad \forall l \in [m]. \quad (3.3)$$

Throughout the training process, we keep the parameter b fixed while updating W . For notational simplicity, we write $u_{W,b}(s, a)$ as $u_W(s, a)$ and $\phi_{W,b}(s, a)$ as $\phi_W(s, a)$ in the sequel. We denote by \mathbb{E}_{init} the expectation taken with respect to the random initialization in (3.3). For an absolute constant $B > 0$, we define the parameter domain as

$$S_B = \{W \in \mathbb{R}^{md} \mid \|W - W_0\|_2 \leq B\}, \quad (3.4)$$

which is the ball centered at W_0 with the domain radius B .

In the sequel, we consider the following energy-based policy,

$$\pi_\theta(a \mid s) = \frac{\exp(\tau \cdot u_\theta(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\tau \cdot u_\theta(s, a'))}, \quad (3.5)$$

where $\tau \geq 0$ is the inverse temperature parameter and $u_\theta(s, a)$ is the energy function parameterized by the neural network defined in (3.1) with $W = \theta$. In the sequel, we call θ the policy parameter. Meanwhile, we parameterize the reward function $r_\beta(s, a)$ as follows,

$$r_\beta(s, a) = (1 - \gamma)^{-1} \cdot u_\beta(s, a), \quad (3.6)$$

where $u_\beta(s, a)$ is the neural network defined in (3.1) with $W = \beta$ and γ is the discount factor. Here we use the scaling parameter $(1 - \gamma)^{-1}$ to ensure that for any policy π the state-action value function $Q_{r_\beta}^\pi(s, a)$ defined in (2.1) is well approximated by the neural network defined in (3.1). In the sequel, we call β the reward parameter and define the reward function class as

$$\mathcal{R}_\beta = \{r_\beta(s, a) \mid \beta \in S_{B_\beta}\},$$

where S_{B_β} is the parameter domain of β defined in (3.4) with domain radius B_β . To facilitate algorithm design, we establish the following proposition, which calculates the explicit expressions of the gradients $\nabla L(\theta, \beta)$ and the Fisher information $\mathcal{I}(\theta)$. Recall that the Fisher information is defined as

$$\mathcal{I}(\theta) = \mathbb{E}_{(s,a) \sim \nu_{\pi_\theta}} [\nabla_\theta \log \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)^\top]. \quad (3.7)$$

Proposition 3.1 (Gradients and Fisher Information). We call $\iota_\theta(s, a) = \tau^{-1} \cdot \nabla_\theta \log \pi_\theta(a \mid s)$ the temperature-adjusted score function. It holds that

$$\iota_\theta(s, a) = \phi_\theta(s, a) - \mathbb{E}_{a' \sim \pi_\theta(\cdot \mid s)} [\phi_\theta(s, a')]. \quad (3.8)$$

It then holds that

$$\mathcal{I}(\theta) = \tau^2 \cdot \mathbb{E}_{(s,a) \sim \nu_{\pi_\theta}} [\iota_\theta(s, a) \iota_\theta(s, a)^\top], \quad (3.9)$$

$$\nabla_\theta L(\theta, \beta) = -\tau \cdot \mathbb{E}_{(s,a) \sim \nu_{\pi_\theta}} [Q_{r_\beta}^{\pi_\theta}(s, a) \cdot \iota_\theta(s, a)], \quad (3.10)$$

$$\begin{aligned} \nabla_\beta L(\theta, \beta) &= (1 - \gamma)^{-1} \cdot \mathbb{E}_{(s,a) \sim \nu_E} [\phi_\beta(s, a)] \\ &\quad - (1 - \gamma)^{-1} \cdot \mathbb{E}_{(s,a) \sim \nu_{\pi_\theta}} [\phi_\beta(s, a)] \\ &\quad - \lambda \cdot \nabla_\beta \psi(\beta), \end{aligned} \quad (3.11)$$

where $Q_{r_\beta}^{\pi_\theta}(s, a)$ is the state-action value function defined in (2.1) with $\pi = \pi_\theta$ and $r = r_\beta$, ν_{π_θ} is the state-action visitation measure defined in (2.3) with $\pi = \pi_\theta$, and $\mathcal{I}(\theta)$ is the Fisher information defined in (3.7).

Proof. See §C.1 for a detailed proof. \square

Note that the expression of the policy gradient $\nabla_\theta L(\theta, \beta)$ in (3.10) of Proposition 3.1 involves the state-action value function $Q_{r_\beta}^{\pi_\theta}(s, a)$. To this end, we estimate the state-action value function $Q_{r_\beta}^{\pi_\theta}(s, a)$ by $\widehat{Q}_\omega(s, a)$, which is parameterized as follows,

$$\widehat{Q}_\omega(s, a) = u_\omega(s, a). \quad (3.12)$$

Here $u_\omega(s, a)$ is the neural network defined in (3.1) with $W = \omega$. In the sequel, we call ω the value parameter.

3.2. GAIL with Alternating Updates

We employ an actor-critic scheme with alternating updates of the policy and the reward function, which is presented in Algorithm 1. Specifically, we update the policy parameter θ via natural policy gradient and update the reward parameter β via gradient ascent in the actor step, while we estimate the state-action value function $Q_{r_\beta}^\pi(s, a)$ via neural temporal difference (TD) (Cai et al., 2019c) in the critic step.

Actor Step. In the k -th actor step, we update the policy parameter θ and the reward parameter β as follows,

$$\theta_{k+1} = \tau_{k+1}^{-1} \cdot (\tau_k \cdot \theta_k - \eta \cdot \delta_k), \quad (3.13)$$

$$\beta_{k+1} = \text{Proj}_{S_{B_\beta}} \{\beta_k + \eta \cdot \widehat{\nabla}_\beta L(\theta_k, \beta_k)\}, \quad (3.14)$$

where $\tau_{k+1} = \eta + \tau_k$ and

$$\delta_k \in \underset{\delta \in S_{B_\theta}}{\text{argmin}} \|\widehat{\mathcal{I}}(\theta_k) \delta - \tau_k \cdot \widehat{\nabla}_\theta L(\theta_k, \beta_k)\|_2. \quad (3.15)$$

Here $\eta > 0$ is the stepsize, S_{B_θ} and S_{B_β} are the parameter domains of θ and β defined in (3.4) with domain radii B_θ and B_β , respectively, $\text{Proj}_{S_{B_\beta}} : \mathbb{R}^{md} \rightarrow S_{B_\beta}$ is the projection operator, τ_k is the inverse temperature parameter of π_{θ_k} , and $\widehat{\mathcal{I}}(\theta_k)$, $\widehat{\nabla}_\theta L(\theta_k, \beta_k)$, $\widehat{\nabla}_\beta L(\theta_k, \beta_k)$ are the estimators of $\mathcal{I}(\theta_k)$, $\nabla_\theta L(\theta_k, \beta_k)$, $\nabla_\beta L(\theta_k, \beta_k)$, respectively, which are defined in the sequel. In (3.13), we update the policy parameter θ_k along the direction δ_k , which approximates the natural policy gradient $\mathcal{I}(\theta)^{-1} \cdot \nabla_\theta L(\theta, \beta)$, and in (3.15) we update the inverse temperature parameter τ_k to ensure that $\theta_{k+1} \in S_{B_\theta}$. Meanwhile, in (3.14), we update the reward parameter β via (projected) gradient ascent. Following from (3.9) and (3.10) of Proposition 3.1, we construct the estimators of $\mathcal{I}(\theta_k)$ and $\nabla_\theta L(\theta_k, \beta_k)$ as follows,

$$\widehat{\mathcal{I}}(\theta_k) = \frac{\tau_k^2}{N} \sum_{i=1}^N \iota_{\theta_k}(s_i, a_i) \iota_{\theta_k}(s_i, a_i)^\top, \quad (3.16)$$

$$\widehat{\nabla}_\theta L(\theta_k, \beta_k) = -\frac{\tau_k}{N} \sum_{i=1}^N \widehat{Q}_{\omega_k}(s_i, a_i) \cdot \iota_{\theta_k}(s_i, a_i), \quad (3.17)$$

where $\{(s_i, a_i)\}_{i \in [N]}$ is sampled from the state-action visitation measure $\nu_{\pi_{\theta_k}}$ given θ_k with the batch size N , and $\widehat{Q}_{\omega_k}(s, a)$ is the estimator of $Q_{r_{\beta_k}}^{\pi_{\theta_k}}(s, a)$ computed in the critic step. Meanwhile, following from (3.11) of Proposition 3.1, we construct the estimator of $\nabla_\beta L(\theta_k, \beta_k)$ as follows,

$$\begin{aligned} \widehat{\nabla}_\beta L(\theta, \beta) &= \frac{1}{N \cdot (1 - \gamma)} \sum_{i=1}^N [\phi_{\beta_k}(s_i^E, a_i^E) - \phi_{\beta_k}(s_i, a_i)] \\ &\quad - \lambda \cdot \nabla_\beta \psi(\beta_k), \end{aligned} \quad (3.18)$$

where $\{(s_i^E, a_i^E)\}_{i \in [N]}$ is the expert trajectory. For notational simplicity, we write $\pi_k = \pi_{\theta_k}$, $r_k(s, a) = r_{\beta_k}(s, a)$, $d_k =$

d_{π_k} and $\nu_k = \nu_{\pi_k}$ for the k -th step hereafter, where π_θ is the policy, $r_\beta(s, a)$ is the reward function, and d_π, ν_π are the visitation measures defined in (2.3).

Critic Step. Note that the estimator $\widehat{\nabla}_\theta L(\theta, \beta)$ in (3.17) involves the estimator $\widehat{Q}_{\omega_k}(s, a)$ of $Q_{r_k}^{\pi_k}(s, a)$. To this end, we parameterize $\widehat{Q}_\omega(s, a)$ as in (3.12) and adapt neural TD (Cai et al., 2019c), which solves the following minimization problem,

$$\omega_k = \operatorname{argmin}_{\omega \in S_{B_\omega}} \mathbb{E}_{(s,a) \sim \rho_k} [\widehat{Q}_\omega(s, a) - \mathcal{T}_{r_k}^{\pi_k} \widehat{Q}_\omega(s, a)]^2. \quad (3.19)$$

Here S_{B_ω} is the parameter domain with domain radius B_ω , ρ_k is the state-action stationary distribution induced by π_k , and $\mathcal{T}_{r_k}^{\pi_k}$ is the Bellman operator. Note that the Bellman operator \mathcal{T}_r^π is defined as follows,

$$\mathcal{T}_r^\pi Q(s, a) = (1 - \gamma) \cdot r(s, a) + \gamma \cdot \mathbb{E}_\pi [Q(s', a') \mid s, a],$$

where the expectation is taken with respect to $s' \sim P(\cdot \mid s, a)$ and $a' \sim \pi(\cdot \mid s')$. In neural TD, we iteratively update the value parameter ω via

$$\begin{aligned} \delta(j) &= Q_{\omega(j)}(s, a) - r(s, a) - \gamma \cdot Q_{\omega(j)}(s', a'), \\ \omega(j+1) &= \operatorname{Proj}_{S_{B_\omega}} \{ \omega(j) - \alpha \cdot \delta(j) \cdot \nabla_\omega Q_{\omega(j)}(s, a) \}, \end{aligned} \quad (3.20)$$

where $\delta(j)$ is the Bellman residual, $\alpha > 0$ is the stepsize, (s, a) is sampled from the state-action stationary distribution ρ_k , and $s' \sim P(\cdot \mid s, a)$, $a' \sim \pi_k(\cdot \mid s')$ are the subsequent state and action. We defer the detailed discussion of neural TD to §B.

To approximately obtain the compatible function approximation (Sutton et al., 2000; Wang et al., 2019), we share the random initialization among the policy π_θ , the reward function $r_\beta(s, a)$, and the state-action value function $\widehat{Q}_\omega(s, a)$. In other words, we set $\theta_0 = \beta_0 = \omega(0) = W_0$ in our algorithm, where W_0 is the random initialization in (3.3). The output of GAIL is the mixed policy $\bar{\pi}$ (Altman, 1999). Specifically, the mixed policy $\bar{\pi}$ of π_0, \dots, π_{T-1} is executed by randomly selecting a policy π_k for $k \in [0 : T-1]$ with equal probability before time $t = 0$ and exclusively following π_k thereafter. It then holds for any reward function $r(s, a)$ that

$$J(\bar{\pi}; r) = \frac{1}{T} \sum_{k=0}^{T-1} J(\pi_k; r). \quad (3.21)$$

4. Main Results

In this section, we first present the assumptions for our analysis. Then, we establish the global optimality and convergence of Algorithm 1.

Algorithm 1 GAIL

Input: Expert trajectory $\{(s_i^E, a_i^E)\}_{i \in [T_E]}$, number of iterations T , number of iterations T_{TD} of neural TD, stepsize η , stepsize α of neural TD, batch size N , and domain radii $B_\theta, B_\omega, B_\beta$.

- 1: **Initialization.** Initialize $b_l \sim \operatorname{Unif}(\{-1, 1\})$ and $[W_0]_l \sim N(0, I_d/d)$ for any $l \in [m]$ and set $\tau_0 \leftarrow 0$, $\theta_0 \leftarrow W_0$, and $\beta_0 \leftarrow W_0$.
- 2: **for** $k = 0, 1, \dots, T-1$ **do**
- 3: Update value parameter ω_k via Algorithm 2 with π_k , r_k , W_0 , b , T_{TD} , and α as the input.
- 4: Sample $\{(s_i, a_i)\}_{i=1}^N$ from the state-action visitation measure ν_k , and estimate $\widehat{\mathcal{I}}(\theta_k)$, $\widehat{\nabla}_\theta L(\theta_k, \beta_k)$, and $\widehat{\nabla}_\beta L(\theta_k, \beta_k)$ via (3.16), (3.17), and (3.18), respectively.
- 5: Solve $\delta_k \leftarrow \operatorname{argmin}_{\delta \in S_\theta} \|\widehat{\mathcal{I}}(\theta_k) \cdot \delta - \tau_k \cdot \widehat{\nabla}_\theta L(\theta_k, \beta_k)\|_2$ and set $\tau_{k+1} \leftarrow \tau_k + \eta$.
- 6: Update policy parameter θ via $\theta_{k+1} \leftarrow \tau_{k+1}^{-1} \cdot (\tau_k \cdot \theta_k - \eta \cdot \delta_k)$.
- 7: Update reward parameter β via $\beta_{k+1} \leftarrow \operatorname{Proj}_{S_{B_\beta}} \{\beta_k + \eta \cdot \widehat{\nabla}_\beta L(\theta_k, \beta_k)\}$.
- 8: **end for**

Output: Mixed policy $\bar{\pi}$ of π_0, \dots, π_{T-1} .

4.1. Assumptions

We impose the following assumptions on the stationary distributions $\varrho_k \in \mathcal{P}(\mathcal{S})$, $\rho_k \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$ and the visitation measures $d_k \in \mathcal{P}(\mathcal{S})$, $\nu_k \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$.

Assumption 4.1. We assume that the following properties hold.

- (a) Let μ be either ρ_k or ν_k . We assume for an absolute constant $c > 0$ and any $y > 0$ and $w \neq 0$ that

$$\mathbb{E}_{(s,a) \sim \mu} [\mathbb{1}\{|w^\top(s, a)| \leq y\}] \leq \frac{c \cdot y}{\|w\|_2}.$$

- (b) We assume for an absolute constant $C_h > 0$ that

$$\begin{aligned} \max_{k \in \mathbb{N}} \left\{ \left\| \frac{dd_E}{dd_k} \right\|_{2, d_k} + \left\| \frac{d\nu_E}{d\nu_k} \right\|_{2, \nu_k} \right\} &\leq C_h, \\ \max_{k \in \mathbb{N}} \left\{ \left\| \frac{dd_E}{d\varrho_k} \right\|_{2, \varrho_k} + \left\| \frac{d\nu_E}{d\rho_k} \right\|_{2, \rho_k} \right\} &\leq C_h. \end{aligned}$$

Here dd_E/dd_k , $d\nu_E/d\nu_k$, $dd_E/d\varrho_k$, and $d\nu_E/d\rho_k$ are the Radon-Nikodym derivatives.

Assumption 4.1 (a) holds when the probability density functions of ρ_k and ν_k are uniformly upper bounded across k . Assumption 4.1 (b) states that the concentrability coefficients

are uniformly upper bounded across k , which is commonly used in the analysis of RL (Szepesvári & Munos, 2005; Munos & Szepesvári, 2008; Antos et al., 2008; Farahmand et al., 2010; Scherrer et al., 2015; Farahmand et al., 2016; Lazaric et al., 2016).

For notational simplicity, we write $u_0(s, a) = u_{W_0}(s, a)$ and $\phi_0(s, a) = \phi_{W_0}(s, a)$, where $u_{W_0}(s, a)$ is the neural network defined in (3.1) with $W = W_0$, $\phi_{W_0}(s, a)$ is the feature vector defined in (3.2) with $W = W_0$, and W_0 is the random initialization in (3.3). We impose the following assumptions on the neural network $u_0(s, a)$ and the transition kernel P .

Assumption 4.2. We assume that the following properties hold.

- (a) Let $\bar{U} = \sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} |u_0(s, a)|$. We assume for absolute constants $M_0 > 0$ and $v > 0$ and any $t > 2M_0$ that

$$\mathbb{E}_{\text{init}}[\bar{U}^2] \leq M_0^2, \quad \mathbb{P}(\bar{U} > t) \leq \exp(-v \cdot t^2). \quad (4.1)$$

- (b) We assume that the transition kernel P belongs to the following class,

$$\begin{aligned} & \widetilde{\mathcal{M}}_{\infty, B_P} \\ & = \left\{ P(s' | s, a) = \int \vartheta(s, a; w)^\top \varphi(s'; w) dq(w) \mid \right. \\ & \quad \left. \sup_w \left\| \int \varphi(s; w) ds \right\|_2 \leq B_P \right\}. \end{aligned}$$

Here $B_P > 0$ is an absolute constant, q is the probability density function of $N(0, I_d/d)$, and $\vartheta(s, a; w)$ is defined as $\vartheta(s, a; w) = \mathbb{1}\{w^\top(s, a) > 0\} \cdot (s, a)$.

Assumption 4.2 (b) states that the MDP belongs to (a variant of) the class of linear MDPs (Yang & Wang, 2019a;b; Jin et al., 2019; Cai et al., 2019b). However, our class of transition kernels is infinite-dimensional, and thus, captures a rich class of MDPs. To understand Assumption 4.2 (a), recall that we initialize the neural network with $[W_0]_l \sim N(0, I_d/d)$ and $b_l \sim \text{Unif}(\{-1, 1\})$ for any $l \in [m]$. Thus, the neural network $u_0(s, a)$ defined in (3.1) with $W = W_0$ converges to a Gaussian process indexed by $(s, a) \in \mathcal{S} \times \mathcal{A}$ as m goes to infinity. Following from the facts that the maximum of a Gaussian process over a compact index set is sub-Gaussian (van de Geer & Muro, 2014) and that $\mathcal{S} \times \mathcal{A}$ is compact, it is reasonable to assume that $\sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} |u_0(s, a)|$ is sub-Gaussian, which further implies the existence of the absolute constants M_0 and v in (4.1) of Assumption 4.2 (a). Moreover, if we assume that m is even and initialize the parameters W_0, b as follows,

$$\begin{cases} [W_0]_l \stackrel{\text{i.i.d.}}{\sim} N(0, I_d/d), \quad b_l \sim \text{Unif}(\{-1, 1\}), \quad \forall l \leq m/2, \\ [W_0]_l = [W_0]_{l-m/2}, \quad b_l = -b_{l-m/2}, \quad \forall l > m/2, \end{cases} \quad (4.2)$$

we have that $u_0(s, a) = 0$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, which allows us to set $M_0 = 0$ and $v = +\infty$ in Assumption 4.2 (a). Also, it holds that $0 = u_0(s, a) \in \mathcal{R}_\beta$, which implies that $\mathbb{D}_{\mathcal{R}_\beta}(\pi_1, \pi_2) \geq 0$ for any π_1 and π_2 . The proof of our results with the random initialization in (4.2) is identical.

Finally, we impose the following assumption on the regularizer $\psi(\beta)$ and the variances of the estimators $\widehat{\mathcal{I}}(\theta)$, $\widehat{\nabla}_\theta L(\theta, \beta)$, and $\widehat{\nabla}_\beta L(\theta, \beta)$ defined in (3.16), (3.17), and (3.18), respectively.

Assumption 4.3. We assume that the following properties hold.

- (a) We assume for an absolute constant $\sigma > 0$ that

$$\begin{aligned} \mathbb{E}_k \left[\left\| \widehat{\mathcal{I}}(\theta_k) W - \mathbb{E}_k[\widehat{\mathcal{I}}(\theta_k) W] \right\|_2^2 \right] & \leq \tau_k^4 \sigma^2 / N, \\ \forall W \in S_{B_\theta}, \end{aligned} \quad (4.3)$$

$$\mathbb{E}_k \left[\left\| \widehat{\nabla}_\theta L(\theta_k, \beta_k) - \mathbb{E}_k[\widehat{\nabla}_\theta L(\theta_k, \beta_k)] \right\|_2^2 \right] \leq \tau_k^2 \sigma^2 / N, \quad (4.4)$$

$$\mathbb{E}_k \left[\left\| \widehat{\nabla}_\beta L(\theta_k, \beta_k) - \mathbb{E}_k[\widehat{\nabla}_\beta L(\theta_k, \beta_k)] \right\|_2^2 \right] \leq \sigma^2 / N, \quad (4.5)$$

where τ_k is the inverse temperature parameter in (3.5), $N \in \mathbb{N}_+$ is the batch size, and S_{B_θ} is the parameter domain of θ defined in (3.4) with the domain radius B_θ . Here the expectation \mathbb{E}_k is taken with respect to the k -th batch, which is drawn from ν_k given θ_k .

- (b) We assume that the regularizer $\psi(\beta)$ in (2.4) is convex and L_ψ -Lipschitz continuous over the compact parameter domain S_{B_β} .

Assumption 4.3 (a) holds when $\widehat{Q}_{\omega_k}(s_i, a_i) \cdot \iota_{\theta_k}(s_i, a_i)$, $\iota_{\theta_k}(s_i, a_i) \iota_{\theta_k}(s_i, a_i)^\top$, and $\phi_{\beta_k}(s_i, a_i)$ have uniformly upper bounded variances across $i \in [m]$ and k , and the Markov chain that generates $\{(s_i, a_i)\}_{i \in [N]}$ mixes sufficiently fast (Zhang et al., 2019a). Similar assumptions are also used in the analysis of policy optimization (Xu et al., 2019a;b). Also, Assumption 4.3 (b) holds for most commonly used regularizers (Ho & Ermon, 2016).

4.2. Global Optimality and Convergence

In this section, we establish the global optimality and convergence of Algorithm 1. The following proposition adapted from (Cai et al., 2019c) characterizes the global optimality and convergence of neural TD, which is presented in Algorithm 2.

Proposition 4.4 (Global Optimality and Convergence of Neural TD). In Algorithm 2, we set $T_{\text{TD}} = \Omega(m)$, $\alpha =$

$\min\{(1-\gamma)/8, m^{-1/2}\}$, and $B_\omega = c \cdot (B_\beta + B_P \cdot (M_0 + B_\beta))$ for an absolute constant $c > 0$. Let π_k, r_k be the input and ω_k be the output of Algorithm 2. Under Assumptions 4.1 and 4.2, it holds for an absolute constant $C_v > 0$ that

$$\begin{aligned} \mathbb{E}_{\text{init}} \left[\left\| Q_{\omega_k}(s, a) - Q_{r_k}^{\pi_k}(s, a) \right\|_{2, \rho_k}^2 \right] & \quad (4.6) \\ = \mathcal{O}(B_\omega^3 \cdot m^{-1/2} + B_\omega^{5/2} \cdot m^{-1/4} + B_\omega^2 \cdot \exp(-C_v \cdot B_\omega^2)). \end{aligned}$$

Here the expectation \mathbb{E}_{init} is taken with respect to the random initialization in (3.3).

Proof. See §B.1 for a detailed proof. \square

The term $B_\omega^2 \cdot \exp(-C_v \cdot B_\omega^2)$ in (4.6) of Proposition 4.4 characterizes the hardness of estimating the state-action value function $Q_{r_k}^{\pi_k}(s, a)$ by the neural network defined in (3.1), which arises because $\|Q_{r_k}^{\pi_k}(s, a)\|_\infty$ is not uniformly upper bounded across k . Note that if we employ the random initialization in (4.2), we have that $C_v = +\infty$. And consequently, such a term vanishes. We are now ready to establish the global optimality and convergence of Algorithm 1.

Theorem 4.5 (Global Optimality and Convergence of GAIL). We set $\eta = 1/\sqrt{T}$ and $B_\omega = c \cdot (B_\beta + B_P \cdot (M_0 + B_\beta))$ for an absolute constant $c > 0$, and $B_\theta = B_\omega$ in Algorithm 1. Let $\bar{\pi}$ be the output of Algorithm 1. Under Assumptions 4.1-4.3, it holds that

$$\begin{aligned} \mathbb{E}[\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \bar{\pi})] & \leq \underbrace{\frac{(1-\gamma)^{-1} \cdot \log |\mathcal{A}| + 13\bar{B}^2 + M_0^2 + 8}{\sqrt{T}}}_{\text{(i)}} \\ & \quad + \underbrace{2\lambda \cdot L_\psi \cdot \bar{B}}_{\text{(ii)}} + \underbrace{\frac{1}{T} \sum_{k=0}^{T-1} \varepsilon_k}_{\text{(iii)}}. \end{aligned} \quad (4.7)$$

Here $\bar{B} = \max\{B_\theta, B_\omega, B_\beta\}$, $\mathbb{D}_{\mathcal{R}_\beta}$ is the \mathcal{R}_β -distance defined in (2.5) with $\mathcal{R}_\beta = \{r_\beta(s, a) \mid \beta \in S_{B_\beta}\}$, the expectation is taken with respect to the random initialization in (3.3) and the T batches, and the error term ε_k satisfies that

$$\begin{aligned} \varepsilon_k & = \underbrace{2\sqrt{2} \cdot C_h \cdot \bar{B} \cdot \sigma \cdot N^{-1/2}}_{\text{(iii.a)}} + \underbrace{\varepsilon_{Q,k}}_{\text{(iii.b)}} \\ & \quad + \underbrace{\mathcal{O}(k \cdot \bar{B}^{3/2} \cdot m^{-1/4} + \bar{B}^{5/4} \cdot m^{-1/8})}_{\text{(iii.c)}}, \end{aligned} \quad (4.8)$$

where C_h is defined in Assumption 4.1, L_ψ and σ are defined in Assumption 4.3, and $\varepsilon_{Q,k} = \mathcal{O}(B_\omega^3 \cdot m^{-1/2} + B_\omega^{5/2} \cdot m^{-1/4} + B_\omega^2 \cdot \exp(-C_v \cdot B_\omega^2))$ is the error induced by neural TD (Algorithm 2).

Proof. See §5 for a detailed proof. \square

The optimality gap in (4.7) of Theorem 4.5 is measured by the expected \mathcal{R}_β -distance $\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \bar{\pi})$ between the expert policy π_E and the learned policy $\bar{\pi}$. Thus, by showing that the optimality gap is upper bounded by $\mathcal{O}(1/\sqrt{T})$, we prove that $\bar{\pi}$ (approximately) outperforms the expert policy π_E in expectation when the number of iterations T goes to infinity. As shown in (4.7) of Theorem 4.5, the optimality gap is upper bounded by the sum of the three terms. Term (i) corresponds to the $1/\sqrt{T}$ rate of convergence of Algorithm 1. Term (ii) corresponds to the bias induced by the regularizer $\lambda \cdot \psi(\beta)$ in the objective function $L(\theta, \beta)$ defined in (2.4). Term (iii) is upper bounded by the sum of the three terms in (4.8) of Theorem 4.5. In detail, term (iii.a) corresponds to the error induced by the variances of $\hat{\mathcal{I}}(\theta)$, $\hat{\nabla}_\theta L(\theta, \beta)$, and $\hat{\nabla}_\beta L(\theta, \beta)$ defined in (4.3), (4.4), and (4.5) of Assumption 4.3, which vanishes as the batch size N in Algorithm 1 goes to infinity. Term (iii.b) is the error of estimating $Q_{r_k}^{\pi_k}(s, a)$ by $\hat{Q}_\omega(s, a)$ using neural TD (Algorithm 2). As shown in Proposition 4.4, the estimation error $\varepsilon_{Q,k}$ vanishes as m and B_ω go to infinity. Term (iii.c) corresponds to the linearization error of the neural network defined in (3.1), which is characterized in Lemma A.2. Following from Theorem 4.5, it holds for $B_\omega = \Omega((C_v^{-1} \cdot \log T)^{1/2})$, $m = \Omega(\bar{B}^{10} \cdot T^6)$, and $N = \Omega(\bar{B}^2 \cdot T \cdot \sigma^2)$ that $\mathbb{E}[\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \bar{\pi})] = \mathcal{O}(T^{-1/2} + \lambda)$, which implies the $1/\sqrt{T}$ rate of convergence of Algorithm 1 (up to the bias induced by the regularizer).

5. Proof of Main Results

In this section, we present the proof of Theorem 4.5, which establishes the global optimality and convergence of Algorithm 1. For notational simplicity, we write $\pi^s(a) = \pi(a \mid s)$ for any policy π , state $s \in \mathcal{S}$, and action $a \in \mathcal{A}$. For any policies π_1, π_2 and distribution μ over \mathcal{S} , we denote the expected Kullback-Leibler (KL) divergence by KL^μ , which is defined as $\text{KL}^\mu(\pi_1 \parallel \pi_2) = \mathbb{E}_{s \sim \mu}[\text{KL}(\pi_1^s \parallel \pi_2^s)]$. For any visitation measures $d_\pi \in \mathcal{P}(\mathcal{S})$ and $\nu_\pi \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$, we denote by \mathbb{E}_{d_π} and \mathbb{E}_{ν_π} the expectations taken with respect to $s \sim d_\pi$ and $(s, a) \sim \nu_\pi$, respectively.

Following from the property of the mixed policy $\bar{\pi}$ in (3.21), we have that

$$\begin{aligned} \mathbb{E}[\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \bar{\pi})] & = \mathbb{E}\left[\max_{\beta' \in S_{B_\beta}} J(\pi_E; r_{\beta'}) - J(\bar{\pi}; r_{\beta'}) \right] \quad (5.1) \\ & = \mathbb{E}\left[\max_{\beta' \in S_{B_\beta}} \frac{1}{T} \sum_{k=0}^{T-1} J(\pi_E; r_{\beta'}) - J(\pi_k; r_{\beta'}) \right]. \end{aligned}$$

We now upper bound the optimality gap in (5.1) by upper bounding the following difference of expected cumulative

rewards,

$$\begin{aligned}
 J(\pi_E; r_{\beta'}) - J(\pi_k; r_{\beta'}) &= \underbrace{J(\pi_E; r_k) - J(\pi_k; r_k)}_{(i)} \\
 &+ \underbrace{L(\theta_k, \beta') - L(\theta_k, \beta_k)}_{(ii)} + \underbrace{\lambda \cdot (\psi(\beta') - \psi(\beta_k))}_{(iii)}, \quad (5.2)
 \end{aligned}$$

where $\beta' \in S_{B_\beta}$ is chosen arbitrarily and $L(\theta, \beta)$ is the objective function defined in (2.4). Following from Assumption 4.3 and the fact that $\beta_k, \beta' \in S_{B_\beta}$, we have that

$$\lambda(\psi(\beta') - \psi(\beta_k)) \leq \lambda L_\psi \|\beta' - \beta_k\|_2 \leq 2\lambda L_\psi B_\beta, \quad (5.3)$$

which upper bounds term (iii) of (5.2). It remains to upper bound terms (i) and (ii) of (5.2), which hinges on the one-point convexity of $J(\pi; r)$ with respect to π and the (approximate) convexity of $L(\theta, \beta)$ with respect to β .

Upper bound of term (i) in (5.2). In what follows, we upper bound term (i) of (5.2). We first introduce the following cost difference lemma (Kakade & Langford, 2002), which corresponds to the one-point convexity of $J(\pi; r)$ with respect to π . Recall that $d_E \in \mathcal{P}(\mathcal{S})$ is the state visitation measure induced by the expert policy π_E .

Lemma 5.1 (Cost Difference Lemma, Lemma 6.1 in (Kakade & Langford, 2002)). For any policy π and reward function $r(s, a)$, it holds that

$$J(\pi_E; r) - J(\pi; r) = (1 - \gamma)^{-1} \mathbb{E}_{d_E} \left[\langle Q_r^\pi(s, \cdot), \pi_E^s - \pi^s \rangle_{\mathcal{A}} \right], \quad (5.4)$$

where γ is the discount factor.

Furthermore, we establish the following lemma, which upper bounds the right-hand side of (5.4) in Lemma 5.1.

Lemma 5.2. Under Assumptions 4.1-4.3, we have that

$$\begin{aligned}
 &\mathbb{E}_{d_E} \left[\langle Q_{r_k}^{\pi_k}(s, \cdot), \pi_E^s - \pi_k^s \rangle_{\mathcal{A}} \right] \\
 &= \eta^{-1} \cdot \text{KL}^{d_E}(\pi_E \| \pi_k) - \eta^{-1} \cdot \text{KL}^{d_E}(\pi_E \| \pi_{k+1}) + \Delta_k^{(i)},
 \end{aligned}$$

where

$$\begin{aligned}
 &\mathbb{E} [|\Delta_k^{(i)}|] \quad (5.5) \\
 &= 2\sqrt{2} \cdot C_h \cdot B_\theta^{1/2} \cdot \sigma^{1/2} \cdot N^{-1/4} + \eta \cdot (M_0^2 + 9B_\theta^2) \\
 &+ \epsilon_{Q,k} + \mathcal{O}(\eta^{-1} \cdot \tau_{k+1} \cdot B_\theta^{3/2} \cdot m^{-1/4} + B_\theta^{5/4} \cdot m^{-1/8}).
 \end{aligned}$$

Here M_0 is defined in Assumption 4.2, σ is defined in Assumption 4.3, N is the batch size in (3.16)-(3.18), and $\epsilon_{Q,k} = \mathcal{O}(B_\omega^3 \cdot m^{-1/2} + B_\omega^{5/2} \cdot m^{-1/4} + B_\omega^2 \cdot \exp(-C_v \cdot B_\omega^2))$ for an absolute constant $C_v > 0$, which depends on the absolute constant v in Assumption 4.2.

Proof. See §C.2 for a detailed proof. \square

Combining Lemmas 5.1 and 5.2, we have that

$$\begin{aligned}
 &J(\pi_E; r_k) - J(\pi_k; r_k) \\
 &\leq \frac{\text{KL}^{d_E}(\pi_E \| \pi_k) - \text{KL}^{d_E}(\pi_E \| \pi_{k+1}) + \eta \cdot \Delta_k^{(i)}}{\eta \cdot (1 - \gamma)}, \quad (5.6)
 \end{aligned}$$

which upper bounds term (i) of (5.2). Here $\Delta_k^{(i)}$ is upper bounded in (5.5) of Lemma 5.2.

Upper bound of term (ii) in (5.2). In what follows, we upper bound term (ii) of (5.2). We first establish the following lemma, which characterizes the (approximate) convexity of $L(\theta, \beta)$ with respect to β .

Lemma 5.3. Under Assumption 4.1, it holds for any $\beta' \in S_{B_\beta}$ that

$$\begin{aligned}
 &\mathbb{E}_{\text{init}} [L(\theta_k, \beta') - L(\theta_k, \beta_k)] \\
 &= \mathbb{E}_{\text{init}} [\nabla_\beta L(\theta_k, \beta_k)^\top (\beta' - \beta_k)] + \mathcal{O}(B_\beta^{3/2} \cdot m^{-1/4}).
 \end{aligned}$$

Proof. See §C.3 for a detailed proof. \square

The term $\mathcal{O}(B_\beta^{3/2} \cdot m^{-1/4})$ in Lemma 5.3 arises from the linearization error of the neural network, which is characterized in Lemma A.2. Based on Lemma 5.3, we establish the following lemma, which upper bounds term (ii) of (5.2).

Lemma 5.4. Under Assumptions 4.1 and 4.3, we have that

$$\begin{aligned}
 &L(\theta_k, \beta') - L(\theta_k, \beta_k) \leq \eta^{-1} \cdot \|\beta_k - \beta'\|_2^2 \\
 &- \eta^{-1} \cdot \|\beta_{k+1} - \beta'\|_2^2 + \Delta_k^{(ii)},
 \end{aligned}$$

where

$$\begin{aligned}
 &\mathbb{E} [|\Delta_k^{(ii)}|] = \eta \cdot \left(2(2(1 - \gamma)^{-1} + \lambda \cdot L_\psi)^2 + \sigma^2 \cdot N^{-1} \right) \\
 &+ 2B_\beta \cdot \sigma \cdot N^{-1/2} + \mathcal{O}(B_\beta^{3/2} \cdot m^{-1/4}). \quad (5.7)
 \end{aligned}$$

Proof. See §C.4 for a detailed proof. \square

By Lemma 5.4, we have that

$$\begin{aligned}
 &L(\theta_k, \beta') - L(\theta_k, \beta_k) \leq \Delta_k^{(ii)} \\
 &+ \eta^{-1} \cdot (\|\beta_k - \beta'\|_2^2 - \|\beta_{k+1} - \beta'\|_2^2), \quad (5.8)
 \end{aligned}$$

which upper bounds term (ii) of (5.2). Here $\Delta_k^{(ii)}$ is upper bounded in (5.7) of Lemma 5.4.

Plugging (5.3), (5.6), and (5.8) into (5.2), we obtain that

$$\begin{aligned}
 &J(\pi_E; r_{\beta'}) - J(\pi_k; r_{\beta'}) \quad (5.9) \\
 &\leq \frac{\text{KL}^{d_E}(\pi_E \| \pi_k) - \text{KL}^{d_E}(\pi_E \| \pi_{k+1})}{\eta \cdot (1 - \gamma)} \\
 &+ \eta^{-1} \|\beta_k - \beta'\|_2^2 - \eta^{-1} \|\beta_{k+1} - \beta'\|_2^2 + 2\lambda \cdot L_\psi \cdot B_\beta + \Delta_k.
 \end{aligned}$$

Here $\Delta_k = \Delta_k^{(i)} + \Delta_k^{(ii)}$, where $\Delta_k^{(i)}$ and $\Delta_k^{(ii)}$ are upper bounded in (5.5) and (5.7) of Lemmas 5.2 and 5.4, respectively. Note that the upper bound of Δ_k does not depend on θ and β . Upon telescoping (5.9) with respect to k , we obtain that

$$\begin{aligned} J(\pi_E; r_{\beta'}) - J(\bar{\pi}; r_{\beta'}) &= \frac{1}{T} \sum_{k=0}^{T-1} [J(\pi_E; r_{\beta'}) - J(\pi_k; r_{\beta'})] \\ &\leq \frac{(1-\gamma)^{-1} \cdot \text{KL}^{d_E}(\pi_E \| \pi_0) + \|\beta_0 - \beta'\|_2^2}{\eta \cdot T} \\ &\quad + 2\lambda \cdot L_\psi \cdot B_\beta + \frac{1}{T} \sum_{k=0}^{T-1} |\Delta_k|. \end{aligned} \quad (5.10)$$

Following from the fact that $\tau_0 = 0$ and the parameterization of π_θ in (3.5), it holds that π_0^s is the uniform distribution over \mathcal{A} for any $s \in \mathcal{S}$. Thus, we have $\text{KL}^{d_E}(\pi_E \| \pi_0) \leq \log |\mathcal{A}|$. Meanwhile, following from the fact that $\beta' \in S_{B_\beta}$, it holds that $\|\beta' - \beta_0\|_2 \leq B_\beta$. Finally, by setting $\eta = T^{-1/2}$, $\tau_k = k \cdot \eta$, and $\bar{B} = \max\{B_\theta, B_\beta\}$ in (5.10), we have that

$$\begin{aligned} \mathbb{E}[\mathbb{D}_{\mathcal{R}_\beta}(\pi_E, \bar{\pi})] &= \mathbb{E}\left[\max_{\beta' \in S_{B_\beta}} J(\pi_E; r_{\beta'}) - J(\bar{\pi}; r_{\beta'})\right] \\ &\leq \frac{(1-\gamma)^{-1} \cdot \log |\mathcal{A}| + 4B_\beta^2}{\eta \cdot T} + 2\lambda \cdot L_\psi \cdot B_\beta \\ &\quad + \frac{\mathbb{E}[\max_{\beta'} \sum_{k=0}^{T-1} |\Delta_k|]}{T} \\ &= \frac{(1-\gamma)^{-1} \cdot \log |\mathcal{A}| + 13\bar{B}^2 + M_0^2 + 8}{\sqrt{T}} + 2\lambda \cdot L_\psi \cdot \bar{B} \\ &\quad + \frac{\sum_{k=0}^{T-1} \varepsilon_k}{T}. \end{aligned}$$

Here ε_k is upper bounded as follows,

$$\begin{aligned} \varepsilon_k &= 2\sqrt{2} \cdot C_h \cdot \bar{B} \cdot \sigma \cdot N^{-1/2} + \varepsilon_{Q,k} \\ &\quad + \mathcal{O}(k \cdot \bar{B}^{3/2} \cdot m^{-1/4} + \bar{B}^{5/4} \cdot m^{-1/8}), \end{aligned}$$

where $\varepsilon_{Q,k} = \mathcal{O}(B_\omega^3 \cdot m^{-1/2} + B_\omega^{5/2} \cdot m^{-1/4} + B_\omega^2 \cdot \exp(-C_v \cdot B_\omega^2))$ for an absolute constant $C_v > 0$. Thus, we complete the proof of Theorem 4.5.

References

- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient methods in Markov decision processes. *arXiv preprint arXiv:1908.00261*, 2019.
- Altman, E. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.
- Anthony, M. and Bartlett, P. L. *Neural network learning: Theoretical foundations*. Cambridge University Press, 2009.
- Antos, A., Szepesvári, C., and Munos, R. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein GAN. *arXiv preprint arXiv:1701.07875*, 2017.
- Bhandari, J. and Russo, D. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.
- Cai, Q., Hong, M., Chen, Y., and Wang, Z. On the global convergence of imitation learning: A case for linear quadratic regulator. *arXiv preprint arXiv:1901.03674*, 2019a.
- Cai, Q., Yang, Z., Jin, C., and Wang, Z. Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*, 2019b.
- Cai, Q., Yang, Z., Lee, J. D., and Wang, Z. Neural temporal-difference learning converges to global optima. *arXiv preprint arXiv:1905.10027*, 2019c.
- Chen, M., Wang, Y., Liu, T., Yang, Z., Li, X., Wang, Z., and Zhao, T. On computation and generalization of generative adversarial imitation learning. *arXiv preprint arXiv:2001.02792*, 2020.
- Farahmand, A.-m., Szepesvári, C., and Munos, R. Error propagation for approximate policy and value iteration. In *Advances in Neural Information Processing Systems*, pp. 568–576, 2010.
- Farahmand, A.-m., Ghavamzadeh, M., Szepesvári, C., and Mannor, S. Regularized policy iteration with nonparametric function spaces. *The Journal of Machine Learning Research*, 17(1):4809–4874, 2016.
- 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, 2016.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Ho, J. and Ermon, S. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, pp. 4565–4573, 2016.
- Hofmann, T., Schölkopf, B., and Smola, A. J. Kernel methods in machine learning. *The Annals of Statistics*, pp. 1171–1220, 2008.

- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019.
- Kakade, S. and Langford, J. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, volume 2, pp. 267–274, 2002.
- Kakade, S. M. A natural policy gradient. In *Advances in Neural Information Processing Systems*, pp. 1531–1538, 2002.
- Kuefler, A., Morton, J., Wheeler, T., and Kochenderfer, M. Imitating driver behavior with generative adversarial networks. In *IEEE Intelligent Vehicles Symposium*, pp. 204–211. IEEE, 2017.
- Lazaric, A., Ghavamzadeh, M., and Munos, R. Analysis of classification-based policy iteration algorithms. *The Journal of Machine Learning Research*, 17(1):583–612, 2016.
- Levine, S. and Koltun, V. Continuous inverse optimal control with locally optimal examples. *arXiv preprint arXiv:1206.4617*, 2012.
- Liu, B., Cai, Q., Yang, Z., and Wang, Z. Neural proximal/trust region policy optimization attains globally optimal policy. *arXiv preprint arXiv:1906.10306*, 2019.
- Merel, J., Tassa, Y., TB, D., Srinivasan, S., Lemmon, J., Wang, Z., Wayne, G., and Heess, N. Learning human behaviors from motion capture by adversarial imitation. *arXiv preprint arXiv:1707.02201*, 2017.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *The Journal of Machine Learning Research*, 9(May):815–857, 2008.
- Nesterov, Y. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Ng, A. Y. and Russell, S. J. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning*, pp. 663–670, 2000.
- Peters, J. and Schaal, S. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- Pomerleau, D. A. Efficient training of artificial neural networks for autonomous navigation. *Neural Computation*, 3(1):88–97, 1991.
- Rafique, H., Liu, M., Lin, Q., and Yang, T. Non-convex min-max optimization: Provable algorithms and applications in machine learning. *arXiv:1810.02060*, 2018.
- Rahimi, A. and Recht, B. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pp. 1177–1184, 2008.
- Rahimi, A. and Recht, B. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. *Advances in Neural Information Processing Systems*, pp. 1313–1320, 2009.
- Ross, S. and Bagnell, D. Efficient reductions for imitation learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 661–668, 2010.
- Ross, S., Gordon, G., and Bagnell, D. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 627–635, 2011.
- Russell, S. Learning agents for uncertain environments. In *Conference on Learning Theory*, pp. 101–103, 1998.
- Scherrer, B., Ghavamzadeh, M., Gabillon, V., Lesner, B., and Geist, M. Approximate modified policy iteration and its application to the game of Tetris. *The Journal of Machine Learning Research*, 16:1629–1676, 2015.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT Press, 2018.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pp. 1057–1063, 2000.
- Syed, U., Bowling, M., and Schapire, R. E. Apprenticeship learning using linear programming. In *International Conference on Machine Learning*, pp. 1032–1039, 2008.
- Szepesvári, C. and Munos, R. Finite time bounds for sampling based fitted value iteration. In *International Conference on Machine Learning*, pp. 880–887. ACM, 2005.
- Tai, L., Zhang, J., Liu, M., and Burgard, W. Socially compliant navigation through raw depth inputs with generative adversarial imitation learning. In *IEEE International Conference on Robotics and Automation*, pp. 1111–1117. IEEE, 2018.
- van de Geer, S. and Muro, A. On higher order isotropy conditions and lower bounds for sparse quadratic forms. *Electronic Journal of Statistics*, 8(2):3031–3061, 2014. doi: 10.1214/15-EJS983.
- Wang, L., Cai, Q., Yang, Z., and Wang, Z. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.

- Xu, P., Gao, F., and Gu, Q. An improved convergence analysis of stochastic variance-reduced policy gradient. *arXiv preprint arXiv:1905.12615*, 2019a.
- Xu, P., Gao, F., and Gu, Q. Sample efficient policy gradient methods with recursive variance reduction. *arXiv preprint arXiv:1909.08610*, 2019b.
- Yang, L. and Wang, M. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, pp. 6995–7004, 2019a.
- Yang, L. F. and Wang, M. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019b.
- Yu, L., Zhang, W., Wang, J., and Yu, Y. SeqGAN: Sequence generative adversarial nets with policy gradient. *arXiv preprint arXiv:1609.05473*, 2016.
- Zhang, K., Koppel, A., Zhu, H., and Başar, T. Global convergence of policy gradient methods to (almost) locally optimal policies. *arXiv preprint arXiv:1906.08383*, 2019a.
- Zhang, K., Yang, Z., and Baar, T. Policy optimization provably converges to Nash equilibria in zero-sum linear quadratic games. *arXiv preprint arXiv:1906.00729*, 2019b.