Error Bounds of Imitating Policies and Environments*

Tian Xu¹, Ziniu Li^{2,3}, Yang Yu^{1,3}

¹National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China ²The Chinese University of Hong Kong, Shenzhen, Shenzhen 518172, China ³Polixir Technologies, Nanjing 210038, China xut@lamda.nju.edu.cn, ziniuli@link.cuhk.edu.cn, yuy@nju.edu.cn

Abstract

Imitation learning trains a policy by mimicking expert demonstrations. Various imitation methods were proposed and empirically evaluated, meanwhile, their theoretical understanding needs further studies. In this paper, we firstly analyze the value gap between the expert policy and imitated policies by two imitation methods, behavioral cloning and generative adversarial imitation. The results support that generative adversarial imitation can reduce the compounding errors compared to behavioral cloning, and thus has a better sample complexity. Noticed that by considering the environment transition model as a dual agent, imitation learning can also be used to learn the environment model. Therefore, based on the bounds of imitating policies, we further analyze the performance of imitating environments. The results show that environment models can be more effectively imitated by generative adversarial imitation than behavioral cloning, suggesting a novel application of adversarial imitation for model-based reinforcement learning. We hope these results could inspire future advances in imitation learning and model-based reinforcement learning.

1 Introduction

Sequential decision-making under uncertainty is challenging due to the stochastic dynamics and delayed feedback [27, 8]. Compared to reinforcement learning (RL) [46, 38] that learns from delayed feedback, imitation learning (IL) [37, 34, 23] learns from expert demonstrations that provide immediate feedback and thus is efficient in obtaining a good policy, which has been demonstrated in playing games [45], robotic control [19], autonomous driving [14], etc.

Imitation learning methods have been designed from various perspectives. For instance, behavioral cloning (BC) [37, 50] learns a policy by directly minimizing the action probability discrepancy with supervised learning; apprenticeship learning (AL) [1, 47] infers a reward function from expert demonstrations by inverse reinforcement learning [34], and subsequently learns a policy by reinforcement learning using the recovered reward function. Recently, Ho and Ermon [23] revealed that AL can be viewed as a state-action occupancy measure matching problem. From this connection, they proposed the algorithm generative adversarial imitation learning (GAIL). In GAIL, a discriminator scores the agent's behaviors based on the similarity to the expert demonstrations, and the agent learns to maximize the scores, resulting in expert-like behaviors.

Many empirical studies of imitation learning have been conducted. It has been observed that, for example, GAIL often achieves better performance than BC [23, 28, 29]. However, the theoretical

^{*}This work is supported by National Key R&D Program of China (2018AAA0101100), NSFC (61876077), and Collaborative Innovation Center of Novel Software Technology and Industrialization. Yang Yu is the corresponding author. This work was done when Ziniu Li was an intern in Polixir Technologies.

explanations behind this observation have not been fully understood. Only until recently, there emerged studies towards understanding the generalization and computation properties of GAIL [13, 55]. In particular, Chen *et al.* [13] studied the generalization ability of the so-called \mathcal{R} -distance given the complete expert trajectories, while Zhang *et al.* [55] focused on the global convergence properties of GAIL under sophisticated neural network approximation assumptions.

In this paper, we present error bounds on the value gap between the expert policy and imitated policies from BC and GAIL, as well as the sample complexity of the methods. The error bounds indicate that the policy value gap is quadratic w.r.t. the horizon for BC, i.e., $1/(1-\gamma)^2$, and cannot be improved in the worst case, which implies large compounding errors [39, 40]. Meanwhile, the policy value gap is only linear w.r.t. the horizon for GAIL, i.e., $1/(1-\gamma)$. Similar to [13], the sample complexity also hints that controlling the complexity of the discriminator set in GAIL could be beneficial to the generalization. But our analysis strikes that a richer discriminator set is still required to reduce the policy value gap. Besides, our results provide theoretical support for the experimental observation that GAIL can generalize well even provided with incomplete trajectories [28].

Moreover, noticed that by regarding the environment transition model as a dual agent, imitation learning can also be applied to learn the transition model [51, 44, 43]. Therefore, based on the analysis of imitating policies, we further analyze the error bounds of imitating environments. The results indicate that the environment model learning through adversarial approaches enjoys a linear policy evaluation error w.r.t. the model-bias, which improves the previous quadratic results [31, 25] and suggests a promising application of GAIL for model-based reinforcement learning.

2 Background

2.1 Markov Decision Process

An infinite-horizon Markov decision process (MDP) [46, 38] is described by a tuple $\mathcal{M}=(\mathcal{S},\mathcal{A},M^*,R,\gamma,d_0)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, and d_0 specifies the initial state distribution. We assume \mathcal{S} and \mathcal{A} are finite. The decision process runs as follows: at time step t, the agent observes a state s_t from the environment and executes an action a_t , then the environment sends a reward $r(s_t,a_t)$ to the agent and transits to a new state s_{t+1} according to $M^*(\cdot|s_t,a_t)$. Without loss of generality, we assume that the reward function is bounded by R_{\max} , i.e., $|r(s,a)| \leq R_{\max}, \forall (s,a) \in \mathcal{S} \times \mathcal{A}$.

A (stationary) policy $\pi(\cdot|s)$ specifies an action distribution conditioned on state s. The quality of policy π is evaluated by its policy value (i.e., cumulative rewards with a discount factor $\gamma \in [0,1)$): $V_{\pi} = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}, a_{t}) | s_{0} \sim d_{0}, a_{t} \sim \pi(\cdot|s_{t}), s_{t+1} \sim M^{*}(\cdot|s_{t}, a_{t}), t = 0, 1, 2, \cdots\right]$. The goal of reinforcement learning is to find an optimal policy π^{*} such that it maximizes the policy value (i.e., $\pi^{*} = \arg\max_{\pi} V_{\pi}$).

Notice that, in an infinite-horizon MDP, the policy value mainly depends on a finite length of the horizon due to the discount factor. The effective planning horizon [38] $1/(1-\gamma)$, i.e., the total discounting weights of rewards, shows how the discount factor γ relates with the effective horizon. We will see that the effective planning horizon plays an important role in error bounds of imitation learning approaches.

To facilitate later analysis, we introduce the discounted stationary state distribution $d_\pi(s)=(1-\gamma)\sum_{t=0}^\infty \gamma^t \Pr(s_t=s;\pi)$, and the discounted stationary state-action distribution $\rho_\pi(s,a)=(1-\gamma)\sum_{t=0}^\infty \gamma^t \Pr(s_t=s,a_t=a;\pi)$. Intuitively, discounted stationary state (state-action) distribution measures the overall "frequency" of visiting a state (state-action). For simplicity, we will omit "discounted stationary" throughout. We add a superscript to value function and state (state-action) distribution, e.g., $V_\pi^{M^*}$, when it is necessary to clarify the MDP.

2.2 Imitation Learning

Imitation learning (IL) [37, 34, 1, 47, 23] trains a policy by learning from expert demonstrations. In contrast to reinforcement learning, imitation learning is provided with action labels from expert policies. We use $\pi_{\rm E}$ to denote the expert policy and Π to denote the candidate policy class throughout this paper. In IL, we are interested in the policy value gap $V_{\pi_{\rm E}}-V_{\pi}$. In the following, we briefly

introduce two popular methods considered in this paper, behavioral cloning (BC) [37] and generative adversarial imitation learning (GAIL) [23].

Behavioral cloning. In the simplest case, BC minimizes the action probability discrepancy with Kullback–Leibler (KL) divergence between the expert's action distribution and the imitating policy's action distribution. It can also be viewed as the maximum likelihood estimation in supervised learning.

$$\min_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\pi_{E}}} \left[D_{\text{KL}} \left(\pi_{E}(\cdot|s), \pi(\cdot|s) \right) \right] := \mathbb{E}_{(s,a) \sim \rho_{\pi_{E}}} \left[\log \left(\frac{\pi_{E}(a|s)}{\pi(a|s)} \right) \right].$$
(1)

Generative adversarial imitation learning. In GAIL, a discriminator D learns to recognize whether a state-action pair comes from the expert trajectories, while a generator π mimics the expert policy via maximizing the rewards given by the discriminator. The optimization problem is defined as:

$$\min_{\pi \in \Pi} \max_{D \in (0,1)^{S \times \mathcal{A}}} \mathbb{E}_{(s,a) \sim \rho_{\pi}} \left[\log \left(D(s,a) \right) \right] + \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}} \left[\log (1 - D(s,a)) \right].$$

When the discriminator achieves its optimum $D^*(s,a) = \rho_{\pi}(s,a)/(\rho_{\pi}(s,a) + \rho_{\pi_E}(s,a))$, we can derive that GAIL is to minimize the state-action distribution discrepancy between the expert policy and the imitating policy with the Jensen-Shannon (JS) divergence (up to a constant):

$$\min_{\pi \in \Pi} D_{JS}(\rho_{\pi_E}, \rho_{\pi}) := \frac{1}{2} \left[D_{KL}(\rho_{\pi}, \frac{\rho_{\pi} + \rho_{\pi_E}}{2}) + D_{KL}(\rho_{\pi_E}, \frac{\rho_{\pi} + \rho_{\pi_E}}{2}) \right]. \tag{2}$$

3 Related Work

In the domain of imitating policies, prior studies [39, 48, 40, 12] considered the finite-horizon setting and revealed that behavioral cloning [37] leads to the compounding errors (i.e., an optimality gap of $\mathcal{O}(T^2)$, where T is the horizon length). DAgger [40] improved this optimality gap to $\mathcal{O}(T)$ at the cost of additional expert queries. Recently, based on generative adversarial network (GAN) [20], generative adversarial imitation learning [23] was proposed and had achieved much empirical success [17, 28, 29, 11]. Though many theoretical results have been established for GAN [5, 54, 3, 26], the theoretical properties of GAIL are not well understood. To the best of our knowledge, only until recently, there emerged studies towards understanding the generalization and computation properties of GAIL [13, 55]. The closest work to ours is [13], where the authors considered the generalization ability of GAIL under a finite-horizon setting with complete expert trajectories. In particular, they analyzed the generalization ability of the proposed \mathcal{R} -distance but they did not provide the bound for policy value gap, which is of interest in practice. On the other hand, the global convergence properties with neural network function approximation were further analyzed in [55].

In addition to BC and GAIL, apprenticeship learning (AL) [1, 47, 49, 24] is a promising candidate for imitation learning. AL infers a reward function from expert demonstrations by inverse reinforcement learning (IRL) [34], and subsequently learns a policy by reinforcement learning using the recovered reward function. In particular, IRL aims to identify a reward function under which the expert policy's performance is optimal. Feature expectation matching (FEM) [1] and game-theoretic apprenticeship learning (GTAL) [47] are two popular AL algorithms with theoretical guarantees under tabular MDP. To obtain an ϵ -optimal policy, FEM and GTAL require expert trajectories of $\mathcal{O}(\frac{k \log k}{(1-\gamma)^2\epsilon^2})$ and $\mathcal{O}(\frac{\log k}{(1-\gamma)^2\epsilon^2})$ respectively, where k is the number of predefined feature functions.

In addition to imitating policies, learning environment transition models can also be treated by imitation learning by considering environment transition model as a dual agent. This connection has been utilized in [44, 43] to model real-world environments and in [51] to reduce the regret regarding model-bias following the idea of DAgger [40], where the model-bias refers to prediction errors when a learned environment model predicts the next state given the current state and current action. Many studies [6, 31, 25] have shown that if the virtual environment is trained with the principle of behavioral cloning (i.e., minimizing one-step transition prediction errors), the learned policy from this learned environment also suffers from the issue of compounding errors regarding model-bias. We are also inspired by the connection between imitation learning and environment-learning, but we focus on applying the distribution matching property of generative adversarial imitation learning to alleviate model-bias. Noticed that in [44, 43], adversarial approaches have been adopted to learn a high fidelity virtual environment, but the reason why such approaches work was unclear. This paper provides a partial answer.

4 Bounds on Imitating Policies

4.1 Imitating Policies with Behavioral Cloning

It is intuitive to understand why behavioral cloning suffers from large compounding errors [48, 40] as that the imitated policy, even with a small training error, may visit a state out of the expert demonstrations, which causes a larger decision error and a transition to further unseen states. Consequently, the policy value gap accumulates along with the planning horizon. The error bound of BC has been established in [48, 40] under a finite-horizon setting, and here we present an extension to the infinite-horizon setting.

Theorem 1. Given an expert policy $\pi_{\rm E}$ and an imitated policy $\pi_{\rm I}$ with $\mathbb{E}_{s \sim d_{\pi_{\rm E}}}[D_{\rm KL}(\pi_{\rm E}(\cdot|s), \pi_{\rm I}(\cdot|s))] \leq \epsilon$ (which can be achieved by BC with objective Eq.(1)), we have that $V_{\pi_{\rm E}} - V_{\pi_{\rm I}} \leq \frac{2\sqrt{2}R_{\rm max}}{(1-\gamma)^2}\sqrt{\epsilon}$.

The proof by the coherent error-propagation analysis can be found in Appendix A. Note that Theorem 1 is under the infinite sample situation. In the finite sample situation, one can further bound the generalization error ϵ in the RHS using classical learning theory (see Corollary 1) and the proof can be found in Appendix A.

Corollary 1. We use $\{(s_{\pi_E}^{(i)}, a_{\pi_E}^{(i)})\}_{i=1}^m$ to denote the expert demonstrations. Suppose that π_E and π_I are deterministic and the provided function class Π satisfies realizability, meaning that $\pi_E \in \Pi$. For policy π_I imitated by BC (see Eq. (1)), $\forall \delta \in (0,1)$, with probability at least $1-\delta$, we have that

$$V_{\pi_{\rm E}} - V_{\pi_{\rm I}} \leq \frac{2R_{\rm max}}{(1-\gamma)^2} \left(\frac{1}{m} \log\left(|\Pi|\right) + \frac{1}{m} \log\left(\frac{1}{\delta}\right)\right).$$

Moreover, we show that the value gap bound in Theorem 1 is tight up to a constant by providing an example shown in Figure 1 (more details can be found in Appendix A.3). Therefore, we conclude that the quadratic dependency on the effective planning horizon, $\mathcal{O}(1/(1-\gamma)^2)$, is inevitable in the worst-case.

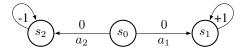


Figure 1: A "hard" deterministic MDP corresponding to Theorem 1. Digits on arrows are corresponding rewards. Initial state is s_0 while s_1 and s_2 are two absorbing states.

4.2 Imitating Policies with GAIL

Different from BC, GAIL [23] is to minimize the state-action distribution discrepancy with JS divergence. The state-action distribution discrepancy captures the temporal structure of Markov decision process, thus it is more favorable in imitation learning. Recent researches [35, 18] showed that besides JS divergence, discrepancy measures based on a general class, f-divergence [30, 16], can be applied to design discriminators. Given two distributions μ and ν , f-divergence is defined as $D_f(\mu,\nu) = \int \mu(x) f(\frac{\mu(x)}{\nu(x)}) dx$, where $f(\cdot)$ is a convex function that satisfies f(1) = 0. Here, we consider GAIL using some common f-divergences listed in Table 1 in Appendix B.1.

Lemma 1. Given an expert policy π_E and an imitated policy π_I with $D_f(\rho_{\pi_I}, \rho_{\pi_E}) \leq \epsilon$ (which can be achieved by GAIL) using the f-divergence in Table 1, we have that $V_{\pi_E} - V_{\pi_I} \leq \mathcal{O}(\frac{1}{1-\gamma}\sqrt{\epsilon})$.

The proof can be found in Appendix B.1. Lemma 1 indicates that the optimality gap of GAIL grows linearly with the effective horizon $1/(1-\gamma)$, multiplied by the square root of the f-divergence error D_f . Compared to Theorem 1, this result indicates that GAIL with the f-divergence could have fewer compounding errors if the objective function is properly optimized. Note that this result does not claim that GAIL is overall better than BC, but can highlight that GAIL has a linear dependency on the planning horizon compared to the quadratic one in BC.

Analyzing the generalization ability of GAIL with function approximation is somewhat more complicated, since GAIL involves a minimax optimization problem. Most of the existing learning theories [32, 42], however, focus on the problems that train one model to minimize the empirical loss, and therefore are hard to be directly applied. In particular, the discriminator in GAIL is often parameterized by certain neural networks, and therefore it may not be optimum within a restrictive function class. In that case, we may view the imitated policy is to minimize the *neural network distance* [5] instead of the ideal *f*-divergence.

Definition 1 (Neural network distance [5]). For a class of neural networks \mathcal{D} , the neural network distance between two (state-action) distributions, μ and ν , is defined as

$$d_{\mathcal{D}}(\mu,\nu) = \sup_{D \in \mathcal{D}} \left\{ \mathbb{E}_{(s,a) \sim \mu}[D(s,a)] - \mathbb{E}_{(s,a) \sim \nu}[D(s,a)] \right\}.$$

Interestingly, it has been shown that the generalization ability of neural network distance is substantially different from the original divergence measure [5, 54] due to the limited representation ability of the discriminator set \mathcal{D} . For instance, JS-divergence may not generalize even with sufficient samples [5]. In the following, we firstly discuss the generalization ability of the neural network distance, based on which we formally give the upper bound of the policy value gap.

To ensure the non-negativity of neural network distance, we assume that the function class \mathcal{D} contains the zero function, i.e., $\exists D \in \mathcal{D}, D(s,a) \equiv 0$. Neural network distance is also known as integral probability metrics (IPM) [33].

As an illustration, f-divergence is connected with neural network distance by its variational representation [54]:

$$d_{f,\mathcal{D}}(\mu,\nu) = \sup_{d \in \mathcal{D}} \left\{ \mathbb{E}_{(s,a) \sim \mu} [D(s,a)] - \mathbb{E}_{(s,a) \sim \nu} [D(s,a)] - \mathbb{E}_{(s,a) \sim \mu} [\phi^*(f(s,a))] \right\},\,$$

where ϕ^* is the (shifted) convex conjugate of f. Thus, considering $\phi^* = 0$ and choosing the activation function of the last layer in the discriminator as the sigmoid function $g(t) = 1/(1 + \exp(-t))$ recovers the original GAIL objective [54]. Again, such defined neural network distance is still different from the original f-divergence because of the limited representation ability of \mathcal{D} . Thereafter, we may consider GAIL is to find a policy $\pi_{\rm I}$ by minimizing $d_{\mathcal{D}}(\rho_{\pi_{\rm E}}, \rho_{\pi_{\rm I}})$.

As another illustration, when $\mathcal D$ is the class of all 1-Lipschitz continuous functions, $d_D(\mu,\nu)$ is the well-known Wasserstein distance [4]. From this viewpoint, we give an instance called Wasserstein GAIL (WGAIL) in Appendix D, where the discriminator in practice is to approximate 1-Lipschitz functions with neural networks. However, note that neural network distance in WGAIL is still distinguished from Wasserstein distance since $\mathcal D$ cannot contain all 1-Lipschitz continuous functions.

In practice, GAIL minimizes the empirical neural network distance $d_{\mathcal{D}}(\hat{\rho}_{\pi_{\rm E}},\hat{\rho}_{\pi})$, where $\hat{\rho}_{\pi_{\rm E}}$ and $\hat{\rho}_{\pi}$ denote the empirical version of population distribution $\rho_{\pi_{\rm E}}$ and ρ_{π} with m samples. To analyze its generalization property, we employ the standard Rademacher complexity technique. The Rademacher random variable σ is defined as $\Pr(\sigma=+1)=\Pr(\sigma=-1)=1/2$. Given a function class \mathcal{F} and a dataset $Z=(z_1,z_2,\cdots,z_m)$ that is i.i.d. dram distribution μ , the empirical Rademacher complexity $\hat{\mathcal{R}}_{\mu}^{(m)}(\mathcal{F})=\mathbb{E}_{\sigma}[\sup_{f\in\mathcal{F}}\frac{1}{m}\sum_{i=1}^m\sigma_i f(z_i)]$ measures the richness of function class \mathcal{F} by the ability to fit random variables [32, 42]. The generalization ability of GAIL is analyzed in [13] under a different definition. They focused on how many trajectories, rather than our focus on state-action pairs, are sufficient to guarantee generalization. Importantly, we further disclose the policy value gap in Theorem 2 based on the neural network distance.

Lemma 2 (Generalization of neural network distance). Consider a discriminator class set \mathcal{D} with Δ -bounded value functions, i.e., $|D(s,a)| \leq \Delta$, for all $(s,a) \in \mathcal{S} \times \mathcal{A}, D \in \mathcal{D}$. Given an expert policy π_E and an imitated policy π_I with $d_{\mathcal{D}}(\hat{\rho}_{\pi_E}, \hat{\rho}_{\pi_I}) - \inf_{\pi \in \Pi} d_{\mathcal{D}}(\hat{\rho}_{\pi_E}, \hat{\rho}_{\pi}) \leq \hat{\epsilon}$, then $\forall \delta \in (0,1)$, with probability at least $1 - \delta$, we have

$$d_{\mathcal{D}}(\rho_{\pi_{\mathsf{E}}}, \rho_{\pi_{\mathsf{I}}}) \leq \underbrace{\inf_{\pi \in \Pi} d_{\mathcal{D}}(\hat{\rho}_{\pi_{\mathsf{E}}}, \hat{\rho}_{\pi})}_{\text{Appr}(\Pi)} + \underbrace{2\hat{\mathcal{R}}_{\rho_{\pi_{\mathsf{E}}}}^{(m)}(\mathcal{D}) + 2\hat{\mathcal{R}}_{\rho_{\pi_{\mathsf{I}}}}^{(m)}(\mathcal{D}) + 12\Delta\sqrt{\frac{\log(2/\delta)}{m}}}_{\text{Estm}(\mathcal{D}, m, \delta)} + \hat{\epsilon}.$$

The proof can be found in Appendix B.3. Here $\mathrm{Appr}(\Pi)$ corresponds to the approximation error induced by the limited policy class Π . $\mathrm{Estm}(\mathcal{D}, m, \delta)$ denotes the estimation error of GAIL regarding

to the complexity of discriminator class and the number of samples. Lemma 2 implies that GAIL generalizes if the complexity of discriminator class \mathcal{D} is properly controlled. Concretely, a simpler discriminator class reduces the estimation error, then tends to reduce the neural network distance. Here we provide an example of neural networks with ReLU activation functions to illustrate this.

Example 1 (Neural Network Discriminator Class). We consider the neural networks with ReLU activation functions $(\sigma_1, \ldots, \sigma_L)$. We use b_s to denote the spectral norm bound and b_n to denote the matrix (2,1) norm bound. The discriminator class consists of L-layer neural networks D_A :

$$\mathcal{D} := \{ D_{\mathbf{A}} : \mathbf{A} = (A_1, \dots, A_L), \|A_i\|_{\sigma} \le b_{\mathbf{s}}, \|A_i^{\top}\|_{2,1} \le b_{\mathbf{n}}, \forall i \in \{1, \dots, L\} \},$$

where $D_{\mathbf{A}}(s,a) = \sigma_L(A_L \cdots \sigma_1(A_1[s^\top,a^\top]^\top))$. Then the spectral normalized complexity $R_{\mathbf{A}}$ of network $D_{\mathbf{A}}$ is $\mathcal{O}(L^{\frac{3}{2}}b_s^{\ L-1}b_n)$ (see [9] for more details). Derived by the Theorem 3.4 in [9] and Lemma 2, with probability at least $1-\delta$, we have

$$d_{\mathcal{D}}(\rho_{\pi_{\mathsf{E}}}, \rho_{\pi_{\mathsf{I}}}) \leq \mathcal{O}\left(\frac{L^{\frac{3}{2}}b_{\mathsf{s}}^{L-1}b_{\mathsf{n}}}{m}\left(1 + \log\left(\frac{m}{L^{\frac{3}{2}}b_{\mathsf{s}}^{L-1}b_{\mathsf{n}}}\right)\right) + \Delta\sqrt{\frac{\log(1/\delta)}{m}}\right) + \inf_{\pi \in \Pi}d_{\mathcal{D}}(\hat{\rho}_{\pi_{\mathsf{E}}}, \hat{\rho}_{\pi}) + \hat{\epsilon}.$$

From a theoretical view, reducing the number of layers L could reduce the spectral normalized complexity $R_{\mathbf{A}}$ and the neural network distance $d_{\mathcal{D}}(\rho_{\pi_{\mathrm{E}}},\rho_{\pi_{\mathrm{I}}})$. However, we did not empirically observe that this operation significantly affects the performance of GAIL. On the other hand, consistent with [28, 29], we find the gradient penalty technique [21] can effectively control the model complexity since this technique gives preference for 1-Lipschitz continuous functions. In this way, the number of candidate functions decreases, and thus the Rademacher complexity of discriminator class is controlled. We also note that the information bottleneck [2] technique helps to control the model complexity and to improve GAIL's performance in practice [36] but the rigorous theoretical explanation is unknown.

However, when the discriminator class is restricted to a set of neural networks with relatively small complexity, it is not safe to conclude that the policy value gap $V_{\pi_E} - V_{\pi_I}$ is small when $d_{\mathcal{D}}(\rho_{\pi_E}, \rho_{\pi_I})$ is small. As an extreme case, if the function class \mathcal{D} only contains constant functions, the neural network distance always equals to zero while the policy value gap could be large. Therefore, we still need a richer discriminator set to distinguish different policies. To substantiate this idea, we introduce the linear span of the discriminator class: $\operatorname{span}(\mathcal{D}) = \{c_0 + \sum_{i=1}^n c_i D_i : c_0, c_i \in \mathbb{R}, D_i \in \mathcal{D}, n \in \mathbb{N}\}$. Furthermore, we assume that $\operatorname{span}(\mathcal{D})$, rather than \mathcal{D} , has enough capacity such that the ground truth reward function r lies in it and define the *compatible coefficient* as:

$$||r||_{\mathcal{D}} = \inf \left\{ \sum_{i=1}^{n} |c_i| : r = \sum_{i=1}^{n} c_i D_i + c_0, \forall n \in \mathbb{N}, c_0, c_i \in \mathbb{R}, D_i \in \mathcal{D} \right\}.$$

Here, $||r||_{\mathcal{D}}$ measures the minimum number of functions in \mathcal{D} required to represent r and $||r||_{\mathcal{D}}$ decreases when the discriminator class becomes richer. Now we present the result on generalization ability of GAIL in the view of policy value gap.

Theorem 2 (GAIL Generalization). *Under the same assumption of Lemma 2 and suppose that the ground truth reward function* r *lies in the linear span of discriminator class, with probability at least* $1 - \delta$, the following inequality holds.

$$V_{\pi_{\mathsf{E}}} - V_{\pi_{\mathsf{I}}} \le \frac{\|r\|_{\mathcal{D}}}{1 - \gamma} (\mathrm{Appr}(\Pi) + \mathrm{Estm}(\mathcal{D}, m, \delta) + \hat{\epsilon}).$$

The proof can be found in Appendix B.4. Theorem 2 discloses that the policy value gap grows linearly with the effective horizon, due to the global structure of state-action distribution matching. This is an advantage of GAIL when only a few expert demonstrations are provided. Moreover, Theorem 2 suggests seeking a trade-off on the complexity of discriminator class: a simpler discriminator class enjoys a smaller estimation error, but could enlarge the compatible coefficient. Finally, Theorem 2 implies the generalization ability holds when provided with some state-action pairs, explaining the phenomenon that GAIL still performs well when only having access to incomplete trajectories [28]. One of the limitations of our result is that we do not deeply consider the approximation ability of the policy class and its computation properties with stochastic policy gradient descent.

5 Bounds on Imitating Environments

The task of environment learning is to recover the transition model of MDP, M_{θ} , from data collected in the real environment M^* , and it is the core of model-based reinforcement learning (MBRL). Although environment learning is typically a separated topic with imitation learning, it has been noticed that learning environment transition model can also be treated by imitation learning [51, 44, 43]. In particular, a transition model takes the state s and action a as input and predicts the next state s', which can be considered as a dual agent, so that the imitation learning can be applied. The learned transition probability $M_{\theta}(s'|s,a)$ is expected to be close to the true probability $M^*(s'|s,a)$. Under the background of MBRL, we assess the quality of the learned transition model M_{θ} by the evaluation error of an arbitrary policy π , i.e., $|V_{\pi}^{M^*} - V_{\pi}^{M_{\theta}}|$, where $V_{\pi}^{M^*}$ is the true value and $V_{\pi}^{M_{\theta}}$ is the value in the learned transition model. Note that we focus on learning the transition model, while assuming the true reward function is always available². For simplicity, we only present the error bounds here and it's feasible to extend our results with the concentration measures to obtain finite sample complexity bounds.

5.1 Imitating Environments with Behavioral Cloning

Similarly, we can directly employ behavioral cloning to minimize the one-step prediction errors when imitating environments, which is formulated as the following optimization problem.

$$\min_{\theta} \mathbb{E}_{(s,a) \sim \rho_{\pi_D}^{M^*}} \left[D_{\mathrm{KL}} \left(M^*(\cdot|s,a), M_{\theta}(\cdot|s,a) \right) \right] := \mathbb{E}_{(s,a) \sim \rho_{\pi_D}^{M^*}, s' \sim M^*(\cdot|s,a)} \left[\log \frac{M^*(s'|s,a)}{M_{\theta}(s'|s,a)} \right],$$

where π_D denotes the data-collecting policy and $\rho_{\pi_D}^{M^*}$ denotes its state-action distribution. We will see that the issue of compounding errors also exists in model-based policy evaluation. Intuitively, if the learned environment cannot capture the transition model globally, the policy evaluation error blows up regarding the model-bias, which degenerates the effectiveness of MBRL. In the following, we formally state this result for self-containing, though similar results have been appeared in [31, 25].

Lemma 3. Given a true MDP with transition model M^* , a data-collecting policy π_D , and a learned transition model M_{θ} with $\mathbb{E}_{(s,a)\sim \rho_{\pi_D}^{M^*}}\left[D_{\mathrm{KL}}\left(M^*(\cdot|s,a),M_{\theta}(\cdot|s,a)\right)\right] \leq \epsilon_m$, for an arbitrary bounded divergence policy π , i.e., $\max_s D_{\mathrm{KL}}\left(\pi(\cdot|s),\pi_D(\cdot|s)\right) \leq \epsilon_{\pi}$, the policy evaluation error is bounded by $|V_{\pi}^{M^*}-V_{\pi}^{M_{\theta}}| \leq \frac{\sqrt{2}R_{\max}}{(1-\gamma)^2}\sqrt{\epsilon_m} + \frac{2\sqrt{2}R_{\max}}{(1-\gamma)^2}\sqrt{\epsilon_{\pi}}$.

Note that the policy evaluation error contains two terms, the inaccuracy of the learned model measured under the state-action distribution of the data-collecting policy π_D , and the policy divergence between π and π_D . We realize that the $1/(1-\gamma)^2$ dependency on the policy divergence ϵ_{π} is inevitable (see also the Theorem 1 in TRPO [41]). Hence, we mainly focus on how to reduce the model-bias term.

5.2 Imitating Environments with GAIL

As shown in Lemma 1, GAIL mitigates the issue of compounding errors via matching the state-action distribution of expert policies. Inspired by this observation, we analyze the error bound of GAIL for environment-learning tasks. Concretely, we train a transition model (also as a policy) that takes state s_t and action a_t as inputs and outputs a distribution over next state s_{t+1} ; at the meantime, we also train a discriminator that learns to recognize whether a state-action-next-state tuple (s_t, a_t, s_{t+1}) comes from the "expert" demonstrations, where the "expert" demonstrations should be explained as the transitions collected by running the data-collecting policy in the true environment. This procedure is summarized in Algorithm 1 in Appendix C.3. It is easy to verify that all the occupancy measure properties of GAIL for imitating policies are reserved by 1) augmenting the action space into the new state space; 2) treating the next state space as the action space when imitating environments. In the following, we show that, by GAIL, the dependency on the effective horizon is only linear in the term of model-bias.

²In the case where the reward function is unknown, we can directly learn the reward function with supervised learning and the corresponding sample complexity is a lower order term compared to the one of learning the transition model [7].

We denote $\mu^{M_{\theta}}$ as the state-action-next-state distribution of the data-collecting policy π_D in the learned transition model, i.e., $\mu^{M_{\theta}}(s,a,s') = M_{\theta}(s'|s,a)\rho^{M_{\theta}}_{\pi_D}(s,a)$; and μ^{M^*} as that in the true transition model. The proof of Theorem 3 can be found in Appendix C.

Theorem 3. Given a true MDP with transition model M^* , a data-collecting policy π_D , and a learned transition model M_{θ} with $D_{\text{JS}}(\mu^{M_{\theta}}, \mu^{M^*}) \leq \epsilon_m$, for an arbitrary bounded divergence policy π , i.e. $\max_s D_{\text{KL}}\big(\pi(\cdot|s), \pi_D(\cdot|s)\big) \leq \epsilon_\pi$, the policy evaluation error is bounded by $|V_{\pi}^{M_{\theta}} - V_{\pi}^{M^*}| \leq \frac{2\sqrt{2}R_{\text{max}}}{1-\gamma} \sqrt{\epsilon_m} + \frac{2\sqrt{2}R_{\text{max}}}{(1-\gamma)^2} \sqrt{\epsilon_\pi}$.

Theorem 3 suggests that recovering the environment transition with a GAIL-style learner can mitigate the model-bias when evaluating policies. We provide the experimental evidence in Section 6.2. Combing this model-imitation technique with all kinds of policy optimization algorithms is an interesting direction that we will explore in the future.

6 Experiments

6.1 Imitating Policies

We evaluate imitation learning methods on three MuJoCo benchmark tasks in OpenAI Gym [10], where the agent aims to mimic locomotion skills. We consider the following approaches: BC [37], DAgger [40], GAIL [23], maximum entropy IRL algorithm AIRL [17] and apprenticeship learning algorithms FEM [1] and GTAL [47]. In particular, FEM and GTAL are based on the improved versions proposed in [24]. Besides GAIL, we also involve WGAIL (see Appendix D) in the comparisons. We run the state-of-the-art algorithm SAC [22] to obtain expert policies. All experiments run with 3 random seeds. Experiment details are given in Appendix E.1.

Study of effective horizon dependency. We firstly compare the methods with different effective planning horizons. All approaches are provided with only 3 expert trajectories, except for DAgger that continues to query expert policies during training. When expert demonstrations are scanty, the impact of the scaling factor $1/(1-\gamma)$ could be significant. The relative performance (i.e., V_π/V_{π_E}) of learned policies under MDPs with different discount factors γ is plotted in Figure 2. Note that the performance of expert policies increases as the planning horizon increases, thus the decrease trends of some curves do not imply the decrease trends of policy values. Exact results and learning curves are given in Appendix E. Though some polices may occasionally outperform experts in short planning horizon settings, we care mostly whether a policy can match the performance of expert policies when the planning horizon increases. We can see that when the planning horizon increases, BC is worse than GAIL, and possibly AIRL, FEM and GTAL. The observation confirms our analysis. Consistent with [23], we also have empirically observed that when BC is provided lots of expert demonstrations, the training error and generalization error could be very small. In that case, the discount scaling factor does not dominate and BC's performance is competitive.

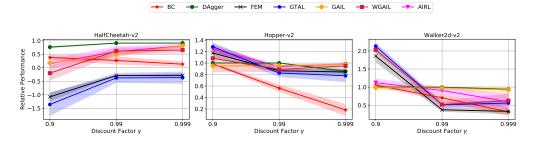


Figure 2: Relative performance of imitated policies under MDPs with different discount factors γ .

Study of generalization ability of GAIL. We then empirically validate the trade-off about the model complexity of the discriminator set in GAIL. We realize that neural networks used by the discriminator are often over-parameterized on MuJoCo tasks and find that carefully using the gradient penalty technique [21] can control the model's complexity to obtain better generalization results. In particular, gradient penalty incurs a quadratic cost function to the gradient norm, which makes the discriminator set a preference to 1-Lipschitz continuous functions. This loss function is multiplied by a coefficient of λ and added to the original objective function. Learning curves of varying λ are given in Figure 3.

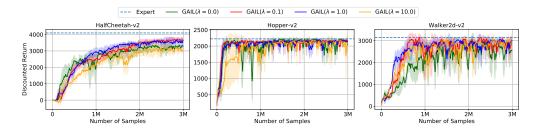


Figure 3: Learning curves of GAIL ($\gamma = 0.999$) with different gradient penalty coefficients λ .

We can see that a moderate λ (e.g., 0.1 or 1.0) yields better performance than a large λ (e.g., 10) or small λ (e.g., 0).

6.2 Imitating Environments

We conduct experiments to verify that generative adversarial learning could mitigate the model-bias for imitating environments in the setting of model-based reinforcement learning. Here we only focus on the comparisons between GAIL and BC for environment-learning. Both methods are provided with 20 trajectories to learn the transition model. Experiment details are given in Appendix E.2. We evaluate the performance by the policy evaluation error $|V_{\pi_D}^{M_\theta} - V_{\pi_D}^{M^*}|$, where π_D denotes the data-collecting policy. As we can see in Figure 4, the policy evaluation errors are smaller on all three environments learned by GAIL. Note that BC tends to over-fit, thus the policy evaluation errors do not decrease on HalfCheetah-v2.

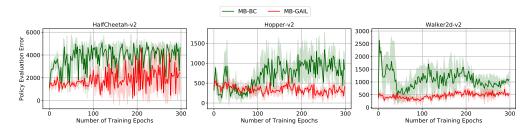


Figure 4: Policy evaluation errors ($\gamma = 0.999$) on environment models trained by BC and GAIL.

7 Conclusion

This paper presents error bounds of BC and GAIL for imitating-policies and imitating-environments in the infinite horizon setting, mainly showing that GAIL can achieve a linear dependency on the effective horizon while BC has a quadratic dependency. The results can enhance our understanding of imitation learning methods.

We would like to highlight that the result of the paper may shed some light for model-based reinforcement learning (MBRL). Previous MBRL methods mostly involve a BC-like transition learning component that can cause a high model-bias. Our analysis suggests that the BC-like transition learner can be replaced by a GAIL-style learner to improve the generalization ability, which also partially addresses the reason that why GAIL-style environment model learning approach in [44, 43] can work well. Learning a useful environment model is an essential way towards sample-efficient reinforcement learning [52], which is not only because the environment model can directly be used for cheap training, but it is also an important support for meta-reinforcement learning (e.g., [53]). We hope this work will inspire future research in this direction.

Our analysis of GAIL focused on the generalization ability of the discriminator and further analysis of the computation and approximation ability of the policy was left for future works.

Broader Impact

This work focuses on the theoretical understanding about imitation learning methods in imitating policies and environments, which does not present any direct societal consequence. This work indicates possible improvement direction for MBRL, which might help reinforcement learning get better used in the real world. There could be some consequence when reinforcement learning is getting abused, such as manipulate information presentation to control people's behaviors.

Acknowledgments and Disclosure of Funding

The authors would like to thank Dr. Weinan Zhang and Dr. Zongzhang Zhang for their helpful comments. This work is supported by National Key R&D Program of China (2018AAA0101100), NSFC (61876077), and Collaborative Innovation Center of Novel Software Technology and Industrialization.

References

- [1] Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the 21st International Conference on Machine Learning (ICML'04)*, pages 1–8, 2004.
- [2] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *Proceedings of the 5th International Conference on Learning Representations (ICLR'17)*, 2017.
- [3] Martín Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. In *Proceedings of the 5th International Conference on Learning Representations* (*ICLR'17*), 2017.
- [4] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning*, (ICML'17), pages 214–223, 2017.
- [5] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (GANs). In *Proceedings of the 34th International Conference on Machine Learning (ICML'17)*, pages 224–232, 2017.
- [6] Kavosh Asadi, Dipendra Misra, and Michael L. Littman. Lipschitz continuity in model-based reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning (ICML'18)*, pages 264–273, 2018.
- [7] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J. Kappen. Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91(3):325–349, 2013.
- [8] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML'17)*, pages 263–272, 2017.
- [9] Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems 30 (NeurIPS'17)*, pages 6240–6249, 2017.
- [10] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. *arXiv*, 1606.01540, 2016.
- [11] Xin-Qiang Cai, Yao-Xiang Ding, Yuan Jiang, and Zhi-Hua Zhou. Imitation learning from pixel-level demonstrations by hashreward. *arXiv*, 1909.03773, 2020.
- [12] Kai-Wei Chang, Akshay Krishnamurthy, Alekh Agarwal, Hal Daumé III, and John Langford. Learning to search better than your teacher. In *Proceedings of the 32nd International Conference on Machine Learning (ICML'15)*, pages 2058–2066, 2015.

- [13] Minshuo Chen, Yizhou Wang, Tianyi Liu, Zhuoran Yang, Xingguo Li, Zhaoran Wang, and Tuo Zhao. On computation and generalization of generative adversarial imitation learning. In *Proceedings of the 8th International Conference on Learning Representations (ICLR'20)*, 2020.
- [14] Felipe Codevilla, Matthias Miiller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *Proceedings of the 2018 IEEE International Conference on Robotics and Automation (ICRA'18)*, pages 1–9, 2018.
- [15] Imre Csiszar and János Körner. *Information theory: coding theorems for discrete memoryless systems*. Cambridge University Press, 2011.
- [16] Imre Csiszár and Paul C. Shields. Information theory and statistics: A tutorial. *Foundations and Trends in Communications and Information Theory*, 1(4), 2004.
- [17] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adverserial inverse reinforcement learning. In *Proceedings of the 6th International Conference on Learning Representations (ICLR'18)*, 2018.
- [18] Seyed Kamyar Seyed Ghasemipour, Richard S. Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. In *Proceedings of the 3rd Conference on Robot Learning (CoRL'19)*, 2019.
- [19] Alessandro Giusti, Jerome Guzzi, Dan C. Ciresan, Fang-Lin He, Juan P. Rodriguez, Flavio Fontana, Matthias Faessler, Christian Forster, Jürgen Schmidhuber, Gianni Di Caro, Davide Scaramuzza, and Luca Maria Gambardella. A machine learning approach to visual perception of forest trails for mobile robots. *IEEE Robotics and Automation Letters*, 1(2):661–667, 2016.
- [20] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems* 27 (NeurIPS'14), pages 2672–2680, 2014.
- [21] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein GANs. In *Advances in Neural Information Processing Systems* 30 (NeurIPS'17), pages 5767–5777, 2017.
- [22] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning (ICML'18)*, pages 1856–1865, 2018.
- [23] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems* 29 (NeurIPS'16), pages 4565–4573, 2016.
- [24] Jonathan Ho, Jayesh K. Gupta, and Stefano Ermon. Model-free imitation learning with policy optimization. In *Proceedings of the 33nd International Conference on Machine Learning (ICML'16)*, pages 2760–2769, 2016.
- [25] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems 32* (*NeurIPS'19*), pages 12498–12509, 2019.
- [26] Haoming Jiang, Zhehui Chen, Minshuo Chen, Feng Liu, Dingding Wang, and Tuo Zhao. On computation and generalization of generative adversarial networks under spectrum control. In Proceedings of the 7th International Conference on Learning Representations (ICLR'19), 2019.
- [27] Michael J. Kearns and Satinder P. Singh. Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 49(2-3):209–232, 2002.
- [28] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning. In *Proceedings of the 7th International Conference on Learning Representations (ICLR'19)*, 2019.
- [29] Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution matching. In *Proceedings of the 8th International Conference on Learning Representations (ICLR*'20), 2020.

- [30] F. Liese and Igor Vajda. On divergences and informations in statistics and information theory. *IEEE Transaction on Information Theory*, 52(10):4394–4412, 2006.
- [31] Yuping Luo, Huazhe Xu, Yuanzhi Li, Yuandong Tian, Trevor Darrell, and Tengyu Ma. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. In *Proceedings of the 7th International Conference on Learning Representations (ICLR'19)*, 2019.
- [32] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. Foundations of Machine Learning. The MIT Press, 2012.
- [33] Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- [34] Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. In Proceedings of the 17th International Conference on Machine Learning (ICML'00), pages 663– 670, 2000.
- [35] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In Advances in Neural Information Processing Systems 29 (NeurIPS'16), pages 271–279, 2016.
- [36] Xue Bin Peng, Angjoo Kanazawa, Sam Toyer, Pieter Abbeel, and Sergey Levine. Variational discriminator bottleneck: Improving imitation learning, inverse RL, and GANs by constraining information flow. In *Proceedings of the 7th International Conference on Learning Representations* (ICLR'19), 2019.
- [37] Dean Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural Computation*, 3(1):88–97, 1991.
- [38] Martin L Puterman. Markov Decision Processes: Discrete Stochastic Dynamic Programming. John Wiley & Sons, 2014.
- [39] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In Proceedings of the 13rd International Conference on Artificial Intelligence and Statistics (AISTATS'10), pages 661–668, 2010.
- [40] Stéphane Ross, Geoffrey J. Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS'11)*, pages 627–635, 2011.
- [41] John Schulman, Sergey Levine, Pieter Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML'15)*, pages 1889–1897, 2015.
- [42] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory To Algorithms*. Cambridge University Press, 2014.
- [43] Wenjie Shang, Yang Yu, Qingyang Li, Zhiwei Qin, Yiping Meng, and Jieping Ye. Environment reconstruction with hidden confounders for reinforcement learning based recommendation. In *Proceedings of the 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'19)*, pages 566–576, 2019.
- [44] Jing-Cheng Shi, Yang Yu, Qing Da, Shi-Yong Chen, and Anxiang Zeng. Virtual-taobao: Virtualizing real-world online retail environment for reinforcement learning. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI'19)*, pages 4902–4909, 2019.
- [45] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P. Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [46] Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. MIT Press, 1998.

- [47] Umar Syed and Robert E. Schapire. A game-theoretic approach to apprenticeship learning. In *Advances in Neural Information Processing Systems 20 (NeurIPS'07)*, pages 1449–1456, 2007.
- [48] Umar Syed and Robert E. Schapire. A reduction from apprenticeship learning to classification. In *Advances in Neural Information Processing Systems 23 (NeurIPS'10)*, pages 2253–2261, 2010.
- [49] Umar Syed, Michael H. Bowling, and Robert E. Schapire. Apprenticeship learning using linear programming. In *Proceedings of the 22nd International Conference on Machine Learning (ICML'08)*, pages 1032–1039, 2008.
- [50] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18), pages 4950–4957, 2018.
- [51] Arun Venkatraman, Martial Hebert, and J. Andrew Bagnell. Improving multi-step prediction of learned time series models. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI'15)*, pages 3024–3030, 2015.
- [52] Yang Yu. Towards sample efficient reinforcement learning. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18)*, pages 5739–5743, 2018.
- [53] Chao Zhang, Yang Yu, and Zhi-Hua Zhou. Learning environmental calibration actions for policy self-evolution. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18)*, pages 3061–3067, 2018.
- [54] Pengchuan Zhang, Qiang Liu, Dengyong Zhou, Tao Xu, and Xiaodong He. On the discrimination-generalization tradeoff in GANs. In *Proceedings of the 6th International Conference on Learning Representations (ICLR'18)*, 2018.
- [55] Yufeng Zhang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Generative adversarial imitation learning with neural networks: Global optimality and convergence rate. *arXiv*, 2003.03709, 2020.

A Analysis of Imitating-policies with BC

Here, we present an error propagation analysis to derive the compounding errors of BC under the setting of infinite-horizon MDP. Our derivation is based on the framework of error-propagation (see Figure 5), which illustrates the cause of compounding errors. Note that the error-propagation framework focuses on the absolute value of policy value gap $|V_{\pi} - V_{\pi_{\rm E}}|$, and the one side bound $V_{\pi} - V_{\pi_{\rm E}}$ can be easily derived from it.

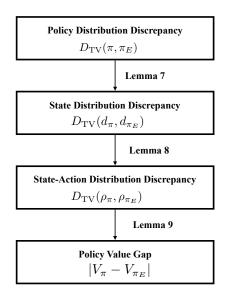


Figure 5: Error propagation of behavioral cloning.

A.1 Error-propagation Analysis

We firstly introduce the following Lemma, which tells that how much state distribution discrepancy grows based on the policy distribution discrepancy.

Lemma 4. For two policies π and π_E , we have that

$$D_{\mathrm{TV}}(d_{\pi}, d_{\pi_{\mathrm{E}}}) \leq \frac{\gamma}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi_{\mathrm{E}}}} \left[D_{\mathrm{TV}} \big(\pi(\cdot | s), \pi_{\mathrm{E}}(\cdot | s) \big) \right].$$

Proof. The proof is based on the permutation theory presented in [41]. First, we show that

$$d_{\pi} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} \Pr(s_{t} = s | \pi, d_{0})$$

= $(1 - \gamma)(I - \gamma P_{\pi})^{-1} d_{0}$,

where $P_{\pi}(s'|s) = \sum_{a \in A} M^*(s'|s,a)\pi(a|s)$. Then we obtain that

$$d_{\pi} - d_{\pi_{\rm E}} = (1 - \gamma)[(I - \gamma P_{\pi})^{-1} - (I - \gamma P_{\pi_{\rm E}})^{-1}] d_0$$

= $(1 - \gamma)(M_{\pi} - M_{\pi_{\rm E}}) d_0,$ (3)

where $M_\pi=(I-\gamma P_\pi)^{-1}$ and $M_{\pi_{\rm E}}=(I-\gamma P_{\pi_{\rm E}})^{-1}$. For the term $M_\pi-M_{\pi_{\rm E}}$, we obtain that

$$M_{\pi} - M_{\pi_{\rm E}} = M_{\pi} \left(M_{\pi_{\rm E}}^{-1} - M_{\pi}^{-1} \right) M_{\pi_{\rm E}}$$

= $\gamma M_{\pi} (P_{\pi} - P_{\pi_{\rm E}}) M_{\pi_{\rm E}}$. (4)

Combining Eq. (3) with Eq. (4), we have

$$d_{\pi} - d_{\pi_{\rm E}} = (1 - \gamma)\gamma M_{\pi} (P_{\pi} - P_{\pi_{\rm E}}) M_{\pi_{\rm E}} d_0$$

= $\gamma M_{\pi} (P_{\pi} - P_{\pi_{\rm E}}) d_{\pi_{\rm E}}$.

Therefore, we obtain that

$$D_{\text{TV}}(d_{\pi}, d_{\pi_{\text{E}}}) = \frac{\gamma}{2} \| M_{\pi} (P_{\pi} - P_{\pi_{\text{E}}}) d_{\pi_{\text{E}}} \|_{1}$$

$$\leq \frac{\gamma}{2} \| M_{\pi} \|_{1} \| (P_{\pi} - P_{\pi_{\text{E}}}) d_{\pi_{\text{E}}} \|_{1}.$$
(5)

We can show that M_{π} is bounded:

$$||M_{\pi}||_{1} = ||\sum_{t=0}^{\infty} \gamma^{t} P_{\pi}^{t}||_{1} \le \sum_{t=0}^{\infty} \gamma^{t} ||P_{\pi}||_{1}^{t} \le \sum_{t=0}^{\infty} \gamma^{t} = \frac{1}{1-\gamma}.$$

Consequently, we show that $\|(P_\pi-P_{\pi_{\rm E}})d_{\pi_{\rm E}}\|_1$ is also bounded,

$$\begin{split} \|(P_{\pi} - P_{\pi_{E}})d_{\pi_{E}}\|_{1} &\leq \sum_{s,s'} |P_{\pi}(s'|s) - P_{\pi_{E}}(s'|s)| \, d_{\pi_{E}}(s) \\ &= \sum_{s,s'} \left| \sum_{a} M^{*}(s'|s,a) \left(\pi(a|s) - \pi_{E}(a|s) \right) \right| \, d_{\pi_{E}}(s) \\ &\leq \sum_{(s,a),s'} M^{*}(s'|s,a) \big| \pi(a|s) - \pi_{E}(a|s) \big| \, d_{\pi_{E}}(s) \\ &= \sum_{s} d_{\pi_{E}}(s) \sum_{a} \big| \pi(a|s) - \pi_{E}(a|s) \big| \\ &= 2\mathbb{E}_{s \sim d_{\pi_{E}}} [D_{\text{TV}} \left(\pi_{E}(\cdot|s), \pi(\cdot|s) \right)]. \end{split}$$

Combining Eq. (5) with the above two inequalities completes the proof.

Next, we further bound the state-action distribution discrepancy based on the policy discrepancy.

Lemma 5. For any two policies π and π_E , we have that

$$D_{\mathrm{TV}}(\rho_{\pi}, \rho_{\pi_{\mathrm{E}}}) \leq \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi_{\mathrm{E}}}} \left[D_{\mathrm{TV}} \left(\pi(\cdot | s), \pi_{\mathrm{E}}(\cdot | s) \right) \right].$$

Proof. Note that the relationship $\rho_{\pi}(s,a) = \pi(a|s)d_{\pi}(s)$ for any policy π , we have

$$\begin{split} &D_{\text{TV}}(\rho_{\pi}, \rho_{\pi_{\text{E}}}) \\ &= \frac{1}{2} \sum_{(s,a)} \left| \left[\pi_{\text{E}}(a|s) - \pi(a|s) \right] d_{\pi_{\text{E}}}(s) + \left[d_{\pi_{\text{E}}}(s) - d_{\pi}(s) \right] \pi(a|s) \right| \\ &\leq \frac{1}{2} \sum_{(s,a)} \left| \pi_{\text{E}}(a|s) - \pi(a|s) \right| d_{\pi_{\text{E}}}(s) + \frac{1}{2} \sum_{(s,a)} \pi(a|s) \left| d_{\pi_{\text{E}}}(s) - d_{\pi}(s) \right| \\ &= \mathbb{E}_{s \sim d_{\pi_{\text{E}}}} [D_{\text{TV}} \left(\pi(\cdot|s), \pi_{\text{E}}(\cdot|s) \right)] + D_{\text{TV}} (d_{\pi}, d_{\pi_{\text{E}}}) \\ &\leq \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi_{\text{E}}}} \left[D_{\text{TV}} \left(\pi(\cdot|s), \pi_{\text{E}}(\cdot|s) \right) \right], \end{split}$$

where the last inequality follows Lemma 4.

Finally, we bound the policy value gap (i.e., the difference between value of learned policy π and the expert policy π_E) based on the state-action distribution discrepancy.

Lemma 6. For any two policies π and π_E , we have that

$$|V_{\pi} - V_{\pi_{\mathsf{E}}}| \le \frac{2R_{\max}}{1 - \gamma} D_{\mathsf{TV}}\left(\rho_{\pi}, \rho_{\pi_{\mathsf{E}}}\right).$$

Proof. It is a well-known fact that for any policy π , its policy value can be reformulated as $V^{\pi} = \frac{1}{1-\gamma}\mathbb{E}_{(s,a)\sim\rho_{\pi}}[r(s,a)]$ [38]. Based on this observation, we derive that

$$\begin{split} |V_{\pi} - V_{\pi_{\mathsf{E}}}| &= \left| \frac{1}{1 - \gamma} \mathbb{E}_{(s, a) \sim \rho_{\pi}}[r(s, a)] - \frac{1}{1 - \gamma} \mathbb{E}_{(s, a) \sim \rho_{\pi_{\mathsf{E}}}}[r(s, a)] \right| \\ &\leq \frac{1}{1 - \gamma} \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left| \left(\rho_{\pi}(s, a) - \rho_{\pi_{\mathsf{E}}}(s, a) \right) r(s, a) \right| \\ &\leq \frac{2R_{\mathsf{max}}}{1 - \gamma} D_{\mathsf{TV}}(\rho_{\pi}, \rho_{\pi_{\mathsf{E}}}). \end{split}$$

A.2 Proof of Theorem 1

Proof of Theorem 1. Suppose that the imitated policy $\pi_{\rm I}$ optimizes the objective of BC up to an ϵ error, i.e., $\mathbb{E}_{s \sim d_{\pi_{\rm E}}}\left[D_{\rm KL}\left(\pi_{\rm I}(\cdot|s), \pi_{\rm E}(\cdot|s)\right)\right] \leq \epsilon$. Combining Lemma 5 and Lemma 6, we have that, for policy $\pi_{\rm I}$ and $\pi_{\rm E}$,

$$\begin{split} V_{\pi_{\mathsf{E}}} - V_{\pi_{\mathsf{I}}} &\leq \frac{2R_{\max}}{1 - \gamma} D_{\mathsf{TV}}(\rho_{\pi_{\mathsf{I}}}, \rho_{\pi_{\mathsf{E}}}) \\ &\leq \frac{2R_{\max}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d_{\pi_{\mathsf{E}}}} \left[D_{\mathsf{TV}} \big(\pi_{\mathsf{I}}(\cdot|s), \pi_{\mathsf{E}}(\cdot|s) \big) \right]. \end{split}$$

Thanks to Pinsker's inequality [15] that for two arbitrary distributions μ and ν , $D_{\rm TV}(\mu,\nu) \le \sqrt{2D_{\rm KL}(\mu,\nu)}$, we obtain that

$$\begin{split} V_{\pi_{\rm E}} - V_{\pi_{\rm I}} &\leq \frac{2R_{\rm max}}{(1-\gamma)^2} \mathbb{E}_{s \sim d_{\pi_{\rm E}}} \left[\sqrt{2D_{\rm KL} \big(\pi_{\rm I}(\cdot|s), \pi_{\rm E}(\cdot|s)\big)} \right] \\ &\leq \frac{2\sqrt{2}R_{\rm max}}{(1-\gamma)^2} \sqrt{\mathbb{E}_{s \sim d_{\pi_{\rm E}}} \left[D_{\rm KL} \big(\pi_{\rm I}(\cdot|s), \pi_{\rm E}(\cdot|s)\big) \right]} \\ &\leq \frac{2\sqrt{2}R_{\rm max}}{(1-\gamma)^2} \sqrt{\epsilon_{\pi}}, \end{split}$$

where the penultimate inequality follows Jensen's inequality $\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)]$, where $\phi(x) = -\sqrt{x}$.

Based on Theorem 1, we provide a sample complexity analysis of BC using classical learning theory.

Proof of Corollary 1. From Lemma 5 and Lemma 6, we obtain that

$$V_{\pi_{\mathsf{E}}} - V_{\pi_{\mathsf{I}}} \le \frac{2R_{\max}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d_{\pi_{\mathsf{E}}}} \left[D_{\mathsf{TV}} \left(\pi_{\mathsf{I}}(\cdot|s), \pi_{\mathsf{E}}(\cdot|s) \right) \right]. \tag{6}$$

Here we consider that $\pi_{\rm I}$ and $\pi_{\rm E}$ are deterministic policies, thus we obtain that

$$\mathbb{E}_{s \sim d_{\pi_{\mathrm{E}}}}[D_{\mathrm{TV}}(\pi(\cdot|s), \pi_{\mathrm{E}}(\cdot|s))] = \mathbb{E}_{s \sim d_{\pi_{\mathrm{E}}}}[\mathbb{I}(\pi(s) \neq \pi_{\mathrm{E}}(s))],$$

where \mathbb{I} is the indicator function. The policy π_{I} is obtained by solving Eq.(1), thus $\pi_{\mathrm{I}}(s_{\pi_{\mathrm{E}}}^{(i)}) = a_{\pi_{\mathrm{E}}}^{(i)}, \forall i \in \{1,\cdots,m\}$. Since behavioral cloning employs supervised learning to learn a policy, we follow the standard argument in the classical learning theory [32] in the remaining proof. We define the expected risk $L(\pi) = \mathbb{E}_{s \sim d_{\pi_{\mathrm{E}}}}[\mathbb{I}(\pi(s) \neq \pi_{\mathrm{E}}(s))]$ and the empirical risk $L_m(\pi) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(\pi(s_{\pi_{\mathrm{E}}}^{(i)}) \neq a_{\pi_{\mathrm{E}}}^{(i)})$. For a fixed $\epsilon > 0$, we define the bad policy class $\Pi_{\mathrm{B}} = \{\pi \in \Pi : L(\pi) > \epsilon\}$. Then we bound the probability of policy π_{I} belongs to the bad policy class Π_{B} :

$$\Pr(L(\pi_{\mathbf{I}}) > \epsilon) = \Pr(\pi_{\mathbf{I}} \in \Pi_{\mathbf{B}}).$$

Because the empirical risk of π_I equals zero, we get that

$$\Pr(\pi_{\mathrm{I}} \in \Pi_{\mathrm{B}}) \leq \Pr(\exists \pi \in \Pi, L_m(\pi) = 0).$$

For a fixed $\pi \in \Pi$, $\Pr(L_m(\pi) = 0) = (1 - L(\pi))^m \le (1 - \epsilon)^m \le e^{-\epsilon m}$, where the last step follows $1 - a \le e^{-a}$. Then we obtain that

$$\Pr(L(\pi_{\mathrm{I}}) > \epsilon) \le \Pr(\exists \pi \in \Pi, L_m(\pi) = 0) \le \sum_{\pi \in \Pi_{\mathrm{B}}} \Pr(L_m(\pi) = 0) \le |\Pi| e^{-\epsilon m}.$$

Setting the right-hand side to be equal to δ , we get that $L(\pi_{\rm I}) \leq \frac{1}{m} (\log(|\Pi|) + \log(\frac{1}{\delta}))$. Combining it with Eq. (6) completes the proof.

A.3 Tightness of Theorem 1

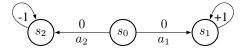


Figure 6: A "hard" deterministic MDP corresponding to Theorem 1. Digits on arrows are corresponding rewards. Initial state is s_0 while s_1 and s_2 are two absorbing states.

Here we validate that the γ -dependence in Theorem 1 is tight by the simple example in Figure 6. Note that the initial state is s_0 and two absorbing states are s_1 and s_2 . That is, the agent always starts with s_0 and takes an action a_1 (a_2); consequently, the system transits into the absorbing state s_1 (s_2). Here we consider a sub-optimal expert policy π_E that chooses a_1 with probability of 0.9 and chooses a_2 with probability of 0.1 at s_0 , meaning that $\pi_E(a_1|s_0)=0.9,\ \pi_E(a_2|s_0)=0.1$ and we can show that the policy value of expert policy π_E is $V_{\pi_E}=\frac{4\gamma}{5(1-\gamma)}$. In addition, we can show that the state distribution of expert policy π_E is $d_{\pi_E}=(d_{\pi_E}(s_0),d_{\pi_E}(s_1),d_{\pi_E}(s_2))=(1-\gamma,\frac{9}{10}\gamma,\frac{1}{10}\gamma)$. Consider a policy obtained by behavioral cloning π_I that chooses a_1 at s_0 with probability of 0.85 and a_2 with probability of 0.15, meaning that $\pi_I(a_1|s_0)=0.85,\pi_I(a_2|s_0)=0.15$. Similarly, we can show that $V_{\pi_I}=\frac{7\gamma}{10(1-\gamma)}$ and the policy value gap $V_{\pi_E}-V_{\pi_I}=\frac{\gamma}{10(1-\gamma)}$. It is easy to verify that the error bound $\mathbb{E}_{s\sim d_{\pi_E}}\left[D_{\mathrm{KL}}\left(\pi_E(\cdot|s),\pi(\cdot|s)\right)\right]$ on the RHS of Eq. (1) is about $0.011(1-\gamma)$ and consequently $V_{\pi_E}-V_{\pi_I}=C\cdot\frac{1}{(1-\gamma)^2}\mathbb{E}_{s\sim d_{\pi_E}}\left[D_{\mathrm{KL}}\left(\pi_E(\cdot|s),\pi(\cdot|s)\right)\right]$, where C is a constant. The equality implies that in the worst case, the quadratic discount complexity is tight in Theorem 1.

B Analysis of Imitating-policies with GAIL

B.1 f-divergence

A large class of divergence measures called f-divergence [30] can be applied to depict the difference between two probability distributions. Given two probability density function μ and ν with respect to a base measure defined on the domain \mathcal{X} , f-divergence is defined as

$$D_f(\mu,\nu) = \int_{\mathcal{X}} \mu(x) f(\frac{\mu(x)}{\nu(x)}) dx,$$

where $f(\cdot)$ is a convex function that satisfies f(1)=0. Different choices of f decides specific measures. When $f(u)=-(u+1)\log(\frac{1+u}{2})+u\log(u)$, f-divergence recovers the JS divergence used in GAIL. Table 1 lists many of the common f-divergences and the f functions to which they correspond (see also [35]). In the following, we provide a proof of Lemma 1. The proof is based on the concentration between different f-divergences.

B.2 Proof of Lemma 1

Proof of Lemma 1. Here we prove that GAIL with f-divergence listed in Table 1 enjoys a linear policy value gap. Derived from Lemma 6, we obtain that

$$V_{\pi_{\rm E}} - V_{\pi} \le \frac{2R_{\rm max}}{1 - \gamma} D_{\rm TV}(\rho_{\pi}, \rho_{\pi_{\rm E}}).$$
 (7)

Table 1: List of f-divergences

Name	$D_f(\mu, u)$	f(u)
Kullback-Leibler	$\int \mu(x) \log(\frac{\mu(x)}{\nu(x)}) dx$ $\int \nu(x) \log(\frac{\nu(x)}{\mu(x)}) dx$ $\int \frac{(\mu(x) - \nu(x))^2}{\mu(x)} dx$	$u\log(u)$
Reverse KL	$\int \nu(x) \log(\frac{\nu(x)}{\mu(x)}) dx$	$-\log(u)$
Pearsion χ^2	$\int \frac{(\mu(x) - \nu(x))^2}{\mu(x)} dx$	$(u-1)^2$
Jensen-Shannon	$\frac{1}{2} \int \mu(x) \log(\frac{2\mu(x)}{\mu(x) + \nu(x)}) + \nu(x) \log(\frac{2\nu(x)}{\mu(x) + \nu(x)}) dx$	$-(u+1)\log(\frac{u+1}{2}) + u\log(u)$
Squared Hellinger	$\int (\sqrt{\mu(x)} - \sqrt{\nu(x)})^2 dx$	$(\sqrt{u}-1)^2$

JS divergence:

In the following, we connect the total variation with the JS divergence based on Pinsker's inequality,

$$D_{\rm JS}(\rho_{\pi_{\rm I}}, \rho_{\pi_{\rm E}}) = \frac{1}{2} \left(D_{\rm KL}(\rho_{\pi_{\rm I}}, \frac{\rho_{\pi_{\rm I}} + \rho_{\pi_{\rm E}}}{2}) + D_{\rm KL}(\rho_{\pi_{\rm E}}, \frac{\rho_{\pi_{\rm I}} + \rho_{\pi_{\rm E}}}{2}) \right)$$

$$\geq D_{\rm TV}^{2}(\rho_{\pi_{\rm I}}, \frac{\rho_{\pi_{\rm I}} + \rho_{\pi_{\rm E}}}{2}) + D_{\rm TV}^{2}(\rho_{\pi_{\rm E}}, \frac{\rho_{\pi_{\rm I}} + \rho_{\pi_{\rm E}}}{2})$$

$$= \frac{1}{2} D_{\rm TV}^{2}(\rho_{\pi_{\rm I}}, \rho_{\pi_{\rm E}}). \tag{8}$$

Combining Eq. (7) with Eq. (8), we get that

$$V_{\pi_{\rm E}} - V_{\pi_{
m I}} \le rac{2\sqrt{2}R_{
m max}}{1-\gamma}\sqrt{D_{
m JS}(
ho_{\pi_{
m I}},
ho_{\pi_{
m E}})}.$$

KL divergence & Reverse KL divergence:

Again, thanks to Pinsker's inequality, we obtain that the policy value gap is bounded by KL divergence and Reverse KL divergence.

$$\begin{split} V_{\pi_{\mathrm{E}}} - V_{\pi_{\mathrm{I}}} &\leq \frac{\sqrt{2}R_{\mathrm{max}}}{1 - \gamma} \sqrt{D_{\mathrm{KL}}(\rho_{\pi_{\mathrm{I}}}, \rho_{\pi_{\mathrm{E}}})}.\\ V_{\pi_{\mathrm{E}}} - V_{\pi_{\mathrm{I}}} &\leq \frac{\sqrt{2}R_{\mathrm{max}}}{1 - \gamma} \sqrt{D_{\mathrm{KL}}(\rho_{\pi_{\mathrm{E}}}, \rho_{\pi_{\mathrm{I}}})}. \end{split}$$

For χ^2 divergence and Squared Hellinger divergence, we can build similar upper bounds of policy value gap.

$$\begin{split} V_{\pi_{\rm E}} - V_{\pi_{\rm I}} &\leq \frac{R_{\rm max}}{1 - \gamma} \sqrt{\chi^2(\rho_{\pi_{\rm I}}, \rho_{\pi_{\rm E}})} \\ V_{\pi_{\rm E}} - V_{\pi_{\rm I}} &\leq \frac{2R_{\rm max}}{1 - \gamma} \sqrt{D_{\rm H}(\rho_{\pi_{\rm I}}, \rho_{\pi_{\rm E}})}. \end{split}$$

In conclusion, for policy $\pi_{\rm I}$ imitated by GAIL with f-divergence listed in Table 1, we have that $V_{\pi_{\rm E}} - V_{\pi_{\rm I}} \leq \mathcal{O}\left(\frac{1}{1-\gamma}\sqrt{D_f\left(\rho_{\pi_{\rm I}},\rho_{\pi_{\rm E}}\right)}\right)$, which finishes the proof.

B.3 Proof of Lemma 2

Proof of Lemma 2. When the policy π_I optimizes the empirical GAIL loss $d_{\mathcal{D}}(\hat{\rho}_{\pi_E}, \hat{\rho}_{\pi})$ up to an ϵ_{opt} error, we have that

$$d_{\mathcal{D}}(\hat{\rho}_{\pi_{\mathsf{E}}}, \hat{\rho}_{\pi_{\mathsf{I}}}) \le \inf_{\pi \in \Pi} d_{\mathcal{D}}(\hat{\rho}_{\pi_{\mathsf{E}}}, \hat{\rho}_{\pi}) + \epsilon_{\mathsf{opt}},\tag{9}$$

where $\hat{\rho}_{\pi_{\rm E}}$ denotes the expert demonstrations with m state-action pairs $\{(s_{\pi_{\rm E}}^{(i)}, a_{\pi_{\rm E}}^{(i)})\}_{i=1}^m$ and $\hat{\rho}_{\pi_{\rm I}}$ is the empirical version of population distribution $\rho_{\pi_{\rm I}}$ with m samples $\{(s_{\pi_{\rm I}}^{(i)}, a_{\pi_{\rm I}}^{(i)})\}_{i=1}^m$ collected by $\pi_{\rm I}$. By standard derivation, we get that

$$d_{\mathcal{D}}(\rho_{\pi_{E}}, \rho_{\pi_{I}}) \le d_{\mathcal{D}}(\rho_{\pi_{E}}, \rho_{\pi_{I}}) - d_{\mathcal{D}}(\hat{\rho}_{\pi_{E}}, \hat{\rho}_{\pi_{I}}) + \inf_{\pi \in \Pi} d_{\mathcal{D}}(\hat{\rho}_{\pi_{E}}, \hat{\rho}_{\pi}) + \epsilon_{\text{opt}}.$$
(10)

According to the definition of neural network distance $d_{\mathcal{D}}(\mu, \nu)$, we prove that $d_{\mathcal{D}}(\rho_{\pi_{E}}, \rho_{\pi_{I}}) - d_{\mathcal{D}}(\hat{\rho}_{\pi_{E}}, \hat{\rho}_{\pi_{I}})$ has an upper bound.

$$\begin{split} & d_{\mathcal{D}}(\rho_{\pi_{\mathsf{E}}},\rho_{\pi_{\mathsf{I}}}) - d_{\mathcal{D}}(\hat{\rho}_{\pi_{\mathsf{E}}},\hat{\rho}_{\pi_{\mathsf{I}}}) \\ & = \sup_{D \in \mathcal{D}} \left[\mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{I}}}}[D(s,a)] \right] - \sup_{D \in \mathcal{D}} \left[\mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{I}}}}[D(s,a)] \right] \\ & \leq \sup_{D \in \mathcal{D}} \left\{ \left[\mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{I}}}}[D(s,a)] \right] - \left[\mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{I}}}}[D(s,a)] \right] \right\} \\ & \leq \sup_{D \in \mathcal{D}} \left[\mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{E}}}}[D(s,a)] \right] + \sup_{D \in \mathcal{D}} \left[\mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{I}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{I}}}}[D(s,a)] \right] \\ & \leq \sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{E}}}}[D(s,a)] \right| + \sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{I}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{\mathsf{I}}}}[D(s,a)] \right|. \end{split}$$

We first show that $\sup_{D\in\mathcal{D}}\left|\mathbb{E}_{(s,a)\sim\rho_{\pi_{\mathrm{E}}}}[D(s,a)]-\mathbb{E}_{(s,a)\sim\hat{\rho}_{\pi_{\mathrm{E}}}}[D(s,a)]\right|$ can be bounded. Note that the assumption that the discriminator set \mathcal{D} consists of bounded functions with Δ , i.e. $\sup_{D\in\mathcal{D}}\|D(s,a)\|_{\infty}\leq 1$

 Δ , $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$. According to McDiarmid 's inequality [32], with probability at least $1 - \frac{\delta}{4}$, the following inequality holds.

$$\sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{E}}} [D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{E}}} [D(s,a)] \right|$$

$$\leq \mathbb{E} \left[\sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{E}}} [D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{E}}} [D(s,a)] \right| \right] + 2\Delta \sqrt{\frac{\log(4/\delta)}{2m}},$$
(11)

where the outer expectation is taken over the random choice of expert demonstrations $\hat{\rho}_{\pi_E}$ with m state-action pairs. According to the Rademacher complexity theory [32], for the first term of Eq. (11) we have that

$$\mathbb{E}\left[\sup_{D\in\mathcal{D}}\left|\mathbb{E}_{(s,a)\sim\rho_{\pi_{E}}}[D(s,a)] - \mathbb{E}_{(s,a)\sim\hat{\rho}_{\pi_{E}}}[D(s,a)]\right|\right] \\
\leq 2\mathbb{E}_{\boldsymbol{\sigma},\rho_{\pi_{E}}}\left[\sup_{D\in\mathcal{D}}\sum_{i=1}^{m}\frac{1}{m}\sigma_{i}D(s^{(i)},a^{(i)})\right] \\
= 2\mathcal{R}_{\rho_{\pi_{E}}}^{(m)}(\mathcal{D}). \tag{12}$$

Based on the connection between Rademacher complexity and empirical Rademacher complexity, we have that with probability at least $1 - \frac{\delta}{4}$, the following inequality holds.

$$\mathcal{R}_{\rho_{\pi_{\mathrm{E}}}}^{(m)}(\mathcal{D}) \le \hat{\mathcal{R}}_{\rho_{\pi_{\mathrm{E}}}}^{(m)}(\mathcal{D}) + 2\Delta\sqrt{\frac{\log(4/\delta)}{2m}},\tag{13}$$

where $\hat{\mathcal{R}}_{\rho_{\pi_{\mathrm{E}}}}^{(m)}(\mathcal{D}) = \mathbb{E}_{\sigma}\left[\sup_{D\in\mathcal{D}}\sum_{i=1}^{m}\frac{1}{m}\sigma_{i}D(s_{\pi_{\mathrm{E}}}^{(i)},a_{\pi_{\mathrm{E}}}^{(i)})\right]$. Combining Eq. (11) with Eq. (13), with probability at least $1-\frac{\delta}{2}$, we have

$$\sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{E}}} [D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_{E}}} [D(s,a)] \right| \le 2\hat{\mathcal{R}}_{\rho_{\pi_{E}}}^{(m)}(\mathcal{D}) + 6\Delta \sqrt{\frac{\log(4/\delta)}{2m}}. \tag{14}$$

By a similar derivation, we obtain that with probability at least $1 - \frac{\delta}{2}$, the following inequality holds.

$$\sup_{D \in \mathcal{D}} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_1}} [D(s,a)] - \mathbb{E}_{(s,a) \sim \hat{\rho}_{\pi_1}} [D(s,a)] \right| \le 2\hat{\mathcal{R}}_{\rho_{\pi_1}}^{(m)}(\mathcal{D}) + 6\Delta \sqrt{\frac{\log(4/\delta)}{2m}}, \tag{15}$$

where $\hat{\mathcal{R}}_{\rho_{\pi_1}}^{(m)}(\mathcal{D}) = \mathbb{E}_{\sigma}\left[\sup_{D\in\mathcal{D}}\sum_{i=1}^{m}\frac{1}{m}\sigma_iD(s_{\pi_1}^{(i)},a_{\pi_1}^{(i)})\right]$. Combining Eq. (10) with Eq. (14) and Eq. (15), we complete the proof.

B.4 Proof of Theorem 2

Proof of Theorem 2. We use the re-formulation of policy value $V_{\pi} = \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim \rho_{\pi}}[r(s,a)]$ and derive that

$$V_{\pi_{\rm E}} - V_{\pi_{\rm I}} \le \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{\rm I}}}[r(s,a)] - \mathbb{E}_{(s,a) \sim \rho_{\pi_{\rm E}}}[r(s,a)] \right|.$$

As we assume that the reward function r lies in the linear span of \mathcal{D} , there exists $n \in \mathbb{N}$, $\{c_i \in \mathbb{R}\}_{i=1}^n$ and $\{D_i \in \mathcal{D}\}_{i=1}^n$, such that $r = c_0 + \sum_{i=1}^n c_i D_i$. Noticed by c_0 will be eliminated by the difference

$$\begin{split} V_{\pi_{\mathsf{E}}} - V_{\pi_{\mathsf{I}}} &\leq \frac{1}{1 - \gamma} \left| \sum_{i=1}^{n} c_{i} \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{I}}}} [D_{i}(s,a)] - \sum_{i=1}^{n} c_{i} \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}} [D_{i}(s,a)] \right| \\ &\leq \frac{1}{1 - \gamma} \sum_{i=1}^{n} \left| c_{i} \right| \left| \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{I}}}} [D_{i}(s,a)] - \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathsf{E}}}} [D_{i}(s,a)] \right| \\ &\leq \frac{1}{1 - \gamma} \left(\sum_{i=1}^{n} |c_{i}| \right) d_{\mathcal{D}}(\rho_{\pi_{\mathsf{I}}}, \rho_{\pi_{\mathsf{E}}}) \\ &\leq \frac{1}{1 - \gamma} \| r \|_{\mathcal{D}} d_{\mathcal{D}}(\rho_{\pi_{\mathsf{I}}}, \rho_{\pi_{\mathsf{E}}}), \end{split}$$

where $||r||_{\mathcal{D}} = \inf\{\sum_{i=1}^{n} |c_i| : r = \sum_{i=1}^{n} c_i D_i + c_0, \forall n \in \mathbb{N}, c_0, c_i \in \mathbb{R}, D_i \in \mathcal{D}\}$. Combining the above inequality with Lemma 2 completes the proof.

\mathbf{C} **Analysis of Imitating-environments**

We first introduce the error bound of policy evaluation without policy divergences, which will be used to prove Lemma 3 later.

Lemma 7. Given an MDP with true transition model M^* , suppose the model error is ϵ_m , i.e., $\mathbb{E}_{(s,a)\sim \rho_{\pi_D}^{M^*}}\left[D_{\mathrm{KL}}\left(M^*(\cdot|s,a),M_{\theta}(\cdot|s,a)\right)\right] \leq \epsilon_m$ (see Eq. (3)), then for the data-collecting policy π_D we have

$$\left| V_{\pi_D}^{M_{\theta}} - V_{\pi_D}^{M^*} \right| \le \frac{\sqrt{2} R_{\text{max}} \gamma}{(1 - \gamma)^2} \sqrt{\epsilon_m}. \tag{16}$$

Proof. The proof is similar to what we have done in Appendix A. First, we show that

$$d_{\pi_D}^{M_{\theta}} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s; \pi_D, M_{\theta}, d_0) = (1 - \gamma)(I - \gamma P_{\theta})^{-1} d_0.$$
 (17)

where $P_{\theta}(s'|s) = \sum_{a \in \mathcal{A}} M_{\theta}(s'|s,a) \pi_D(a|s)$. Following the similar algebraic transformation in Lemma 4, we obtain that

$$d_{\pi_D}^{M_{\theta}} - d_{\pi_D}^{M^*} = \gamma G(P_{\theta} - P^*) d_{\pi_D}^{M^*}$$

 $d_{\pi_D}^{M_\theta}-d_{\pi_D}^{M^*}=\gamma G(P_\theta-P^*)d_{\pi_D}^{M^*},$ where $G_\theta=(I-\gamma P_\theta)^{-1}$ and $G^*=(I-\gamma P^*)^{-1}$. Based on the Cauchy–Schwarz inequality, we have that

$$D_{\text{TV}}(d_{\pi_D}^{M_{\theta}}, d_{\pi_D}^{M^*}) = \frac{\gamma}{2} \|G_{\theta}(P_{\theta} - P^*) d^{M^*}\|_1 \le \frac{\gamma}{2} \|G_{\theta}\|_1 \|(P_{\theta} - P^*) d_{\pi_D}^{M^*}\|_1.$$

We first show that $||G_{\theta}||_1$ is bounded as

$$\|G_{\theta}\|_{1} = \|\sum_{t=0}^{\infty} \gamma^{t} P_{\theta}^{t}\|_{1} \leq \sum_{t=0}^{\infty} \gamma^{t} \|P_{\theta}\|_{1}^{t} \leq \sum_{t=0}^{\infty} \gamma^{t} = \frac{1}{1-\gamma}.$$

We then show that $\|(P_{\theta}-P^*)d_{\pi_D}^{M^*}\|_1$ is bounded,

$$\begin{split} \left\| (P_{\theta} - P^*) d_{\pi_D}^{M^*} \right\|_1 &\leq \sum_{s',s} |P_{\theta}(s'|s) - P^*(s'|s)| d_{\pi_D}^{M^*}(s) \\ &\leq \sum_{s',s,a} |M_{\theta}(s'|s,a) - M^*(s'|s,a)| \pi_D(a|s) d_{\pi_D}^{M^*}(s) \\ &= 2\mathbb{E}_{(s,a) \sim \rho_{\pi_D}^{M^*}} [D_{\text{TV}}(M_{\theta}(\cdot|s,a), M^*(\cdot|s,a))]. \end{split}$$

Thanks to Pinsker's inequality and Jensen's inequality, we can get that

$$D_{\text{TV}}(d_{\pi_D}^{M_{\theta}}, d_{\pi_D}^{M^*}) \le \frac{\sqrt{2}\gamma}{2(1-\gamma)} \sqrt{\epsilon_m}.$$
(18)

From Lemma 6, we obtain that

$$\begin{split} \left| V_{\pi_D}^{M_\theta} - V_{\pi_D}^{M^*} \right| &\leq \frac{R_{\text{max}}}{1 - \gamma} \sum_{(s, a)} \left| \rho_{\pi_D}^{M_\theta}(s, a) - \rho_{\pi_D}^{M^*}(s, a) \right| \\ &\leq \frac{R_{\text{max}}}{1 - \gamma} \sum_{s} \left| d_{\pi_D}^{M_\theta}(s) - d_{\pi_D}^{M^*}(s) \right| \sum_{a} \pi_D(a|s) \\ &\leq \frac{\sqrt{2} R_{\text{max}} \gamma}{(1 - \gamma)^2} \sqrt{\epsilon_m}, \end{split}$$

which concludes the proof.

C.1 Proof of Lemma 3

Proof of Lemma 3. Derived by the triangle inequality, the evaluation error can be decomposed into three parts.

$$|V_{\pi}^{M^*} - V_{\pi}^{M_{\theta}}| \le |V_{\pi}^{M^*} - V_{\pi_{D}}^{M^*}| + |V_{\pi_{D}}^{M^*} - V_{\pi_{D}}^{M_{\theta}}| + |V_{\pi_{D}}^{M_{\theta}} - V_{\pi}^{M_{\theta}}|.$$

For the second term on the RHS, according to Lemma 7, we have

$$|V_{\pi_D}^{M^*} - V_{\pi_D}^{M_\theta}| \le \frac{\sqrt{2}R_{\max}\gamma}{(1-\gamma)^2}\sqrt{\epsilon_m}.$$

For the first term, applying Lemma 5 and Lemma 6, we get that

$$|V_{\pi}^{M^*} - V_{\pi_D}^{M^*}| \leq \frac{2R_{\max}}{1 - \gamma} D_{\text{TV}}(\rho_{\pi}^{M^*}, \rho_{\pi_D}^{M^*})$$

$$\leq \frac{2R_{\max}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d_{\pi_D}^{M^*}} [D_{\text{TV}}(\pi(\cdot|s), \pi_D(\cdot|s))]$$

$$\leq \frac{\sqrt{2}R_{\max}}{(1 - \gamma)^2} \sqrt{\epsilon_{\pi}}.$$

Similar results hold for the third term, meaning that $|V_{\pi_D}^{M_\theta} - V_{\pi}^{M_\theta}| \leq \frac{\sqrt{2}R_{\max}}{(1-\gamma)^2}\sqrt{\epsilon_\pi}$. Combining the above three bounds completes the proof.

C.2 Proof of Theorem 3

Proof of Theorem 3. Due to Lemma 6, we obtain that

$$\begin{split} |V_{\pi}^{M} - V_{\pi}^{M^{*}}| &\leq \frac{2R_{\max}}{1 - \gamma} D_{\text{TV}}(\rho_{\pi}^{M}, \rho_{\pi}^{M^{*}}) \\ &\leq \frac{2R_{\max}}{1 - \gamma} (D_{\text{TV}}(\rho_{\pi}^{M}, \rho_{\pi_{D}}^{M}) + D_{\text{TV}}(\rho_{\pi_{D}}^{M}, \rho_{\pi_{D}}^{M^{*}}) + D_{\text{TV}}(\rho_{\pi_{D}}^{M^{*}}, \rho_{\pi}^{M^{*}})). \end{split}$$

The last inequality follows the triangle inequality. For the term $D_{\text{TV}}(\rho_{\pi_D}^M, \rho_{\pi_D}^{M^*})$, we obtain that

$$D_{\text{TV}}(\rho_{\pi_D}^M, \rho_{\pi_D}^{M^*}) = \frac{1}{2} \sum_{s, a} \left| \sum_{s'} \left(\mu^M(s, a, s') - \mu^{M^*}(s, a, s') \right) \right|$$

$$\leq \frac{1}{2} \sum_{s, a, s'} \left| \mu^M(s, a, s') - \mu^{M^*}(s, a, s') \right|$$

$$= D_{\text{TV}}(\mu^M, \mu^{M^*}).$$

From Eq.(8), we derive that

$$D_{\text{TV}}(\rho_{\pi_D}^M, \rho_{\pi_D}^{M^*}) \leq \sqrt{2D_{\text{JS}}(\rho_{\pi_D}^M, \rho_{\pi_D}^{M^*})} \leq \sqrt{2D_{\text{JS}}(\mu^M, \mu^{M^*})}.$$

Derived by Lemma 5, for the first term $D_{\text{TV}}(\rho_{\pi}^{M}, \rho_{\pi_{D}}^{M})$, we get that

$$\begin{split} D_{\text{TV}}(\rho_{\pi}^{M}, \rho_{\pi_{D}}^{M}) &\leq \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi}^{M}} \left[D_{\text{TV}} \left(\pi(\cdot|s), \pi_{D}(\cdot|s) \right) \right] \\ &\leq \frac{\sqrt{2}}{2(1 - \gamma)} \mathbb{E}_{s \sim d_{\pi}^{M}} \left[\sqrt{D_{\text{KL}} \left(\pi(\cdot|s), \pi_{D}(\cdot|s) \right)} \right]. \\ &\leq \frac{\sqrt{2}}{2(1 - \gamma)} \sqrt{\epsilon_{\pi}}. \end{split}$$

The last two inequalities follow Pinsker's inequality and the definition of ϵ_{π} respectively. Similarly, for the second term $D_{\text{TV}}(\rho_{\pi_D}^{M^*}, \rho_{\pi}^{M^*})$, we have that

$$D_{\text{TV}}(\rho_{\pi_D}^{M^*}, \rho_{\pi}^{M^*}) \le \frac{\sqrt{2}}{2(1-\gamma)} \sqrt{\epsilon_{\pi}}.$$

Combining the above three upper bounds completes the proof.

C.3 Environment-learning with GAIL

Algorithm 1 Environment-learning with GAIL

- 1: **Input:** data-collecting policy π_D , total iterations N, model update iteration N_G , discriminator update iteration N_D .
- 2: Initialize discriminator D, model M_{θ} , and empty dataset \mathcal{B}^* as well as \mathcal{B} .
- 3: $\mathcal{B}^* \leftarrow \text{Collect samples using } \pi_D \text{ in model } M^*.$
- 4: for N iterations do
- 5: **for** N_G iterations **do**
 - $\mathcal{B} \leftarrow \text{Collect samples using } \pi_D \text{ in model } M_{\theta}.$
- 7: Assign rewards to state-action-next-state pairs in \mathcal{B} by discriminator D.
- 8: Update model M_{θ} by maximizing rewards with samples from \mathcal{B} .
- 9: **end for**

6:

- 10: **for** N_D iterations **do**
- 11: Update discriminator D by maximizing the following function:

$$\sum_{(s,a,s')\in\mathcal{B}}[\log(D(s,a,s'))] + \sum_{(s,a,s')\in\mathcal{B}^*}[\log(1-D(s,a,s'))].$$

- 12: end for
- 13: **end for**
- 14: Output: environment model M_{θ} .

The process of applying GAIL to learn the environment transition model is summarized in Algorithm 1.

D Wasserstein GAIL

Similar to Wasserstein GAN (WGAN) [4], we can also introduce Wasserstein distance into GAIL. We call such an algorithm as Wasserstein GAIL (WGAIL for short). Specifically, the discriminator is selected from all 1-Lipschitz function classes by considering the following optimization problem.

$$\max_{D \in ||D||_{\mathrm{Lip}} \leq 1} \mathbb{E}_{(s,a) \sim \rho_{\pi_{\mathrm{E}}}}[D(s,a)] - \mathbb{E}_{(s,a) \sim \rho_{\pi}}[D(s,a)]$$

Due to computation intractability, we cannot compute all 1-Lipschitz functions in practice, and thus we are shifted to its neural network approximation, where D is parameterized by certain neural networks. As our result suggests, this method can still generalize well when its model complexity is controlled. However, ordinary neural networks are often not Lipschitz continuous. To maintain a good approximation to 1-Lipschitz continuous function classes, the gradient penalty technique is

Algorithm 2 Wasserstein GAIL

```
1: Input: Expert demonstrations \mathcal{B}^*, total iterations N, policy update iterations N_G, discriminator
    update iterations N_D.
    Initialize discriminator D, policy \pi, and an empty dataset \mathcal{B}.
 3: for N iterations do
        for N_G iterations do
             \mathcal{B} \leftarrow \text{Collect samples using policy } \pi.
 5:
             Assign scaled rewards to state-action pairs in \mathcal{B} by discriminator D.
 6:
 7:
             Update policy \pi by maximizing rewards with samples from \mathcal{B}.
 8:
        end for
 9:
        for N_D iterations do
10:
             Update discriminator D by maximizing Eq. (19) with samples from \mathcal{B}^* and \mathcal{B}.
        end for
11:
12: end for
13: Output: policy \pi.
```

Table 2: Information about tasks in imitating policies.

Tasks	State Dimension	Action Dimension	Episode Length
HalfCheetah-v2	17	6	1000
Hopper-v2	11	3	1000
Walker2d-v2	17	6	1000

introduced in WGAN [21]. This technique adds a regularization term that employs a quadratic cost to the gradient norm. Hence, denoting (s, a) as z, the loss function for the discriminator in WGAIL is:

$$L(D) = \mathbb{E}_{z \sim \rho_{\pi}} \left[D(z) \right] - \mathbb{E}_{z \sim \rho_{\pi_{\mathbb{F}}}} \left[D(z) \right] + \lambda \mathbb{E}_{z \sim \tilde{\rho}} \left[\left(||\nabla_{z} D(z)|| - 1 \right)^{2} \right], \tag{19}$$

where $\tilde{\rho}$ is a mixing distribution of ρ_{π} and $\rho_{\pi_{E}}$, and λ is a positive regularization coefficient ($\lambda=10$ performs well in practice). Following [24], we also scale reward function (discriminator's output) properly to stabilize training. This is important because the optimization in WGAIL is different from the one in WGAN. Concretely, reinforcement learning algorithms often use the evaluation value rather than the gradient information to perform gradient descent. Without scaling, rewards given by the discriminator often fluctuate, which may lead to an unstable optimization. To tackle this issue, at each iteration, we firstly centralize the given rewards by subtracting the mean and subsequently scale them by dividing the range (the difference between the maximal value the minimal value). The algorithm procedure is outlined in Algorithm 2.

E Experiment Details

E.1 Imitating Policies

We evaluate the considered algorithms on OpenAI Gym [10] benchmark tasks. Information about state dimension, action dimension, and episode length information is listed in Table 2. We run the state-of-the-art algorithm SAC [22] for 1 million samples to obtain expert policies. All imitation learning approaches use 2-layer MLP policy network with 100 hidden sizes and tanh activation function. Except for DAgger that continues to collect new samples and query expert policies (i.e., DAgger collects 1000 samples and gets action labels from expert policies per 5000 iterations), all methods are provided the same 3 expert trajectories with length 1000. Key parameters of BC and DAgger are give in Table 3 and Table 4, respectively. Other methods including GAIL, FEM [1] and GTAL [49] use TRPO [41] to optimize policies, and key parameters are given in Table 5. All experiments run with 3 random seeds (namely, 100, 200 and 300). During the training process, we periodically evaluate the learned policies on true environments with 10 trajectories. Learning curves are given in Figure 7. The final performance of imitated policies and expert policies are listed in Table 6. Please refer to our source code in supplementary materials for other details.

Table 3: Key parameters of Behavioral Cloning.

Parameter	Value		
learning rate	3e-4		
batch size	128		
total number of iters	100k		

Table 4: Key parameters of DAgger.

Parameter	Value
learning rate	3e-4
batch size	128
number of total training iterations	100k
collecting frequency	5k
number of new demonstrations per iteration	1k

Table 5: Key parameters of GAIL, AIRL, FEM and GTAL.

Parameter	Value
number of generator iterations	5
number of discriminator iterations	1
number of rollout samples per iteration	1k
total number of collecting samples	3M
maximal KL divergence	0.01

Table 6: Discounted returns of learned policies. We use \pm to denote the standard deviation

	Tasks	Expert	BC	DAgger	FEM	GTAL	GAIL	WGAIL	AIRL
$\gamma = 0.9$	HalfCheetah-v2 Hopper-v2 Walker2d-v2	$\begin{array}{c} 10.59 \pm 0.00 \\ 10.85 \pm 0.00 \\ 5.31 \pm 0.00 \end{array}$	4.02 ± 0.97 10.69 ± 0.10 5.60 ± 0.20	$\begin{array}{c} 8.06 \pm 0.33 \\ 10.86 \pm 0.02 \\ 5.30 \pm 0.11 \end{array}$	-11.40 ± 1.67 12.75 ± 0.81 9.86 ± 0.87	-14.39 ± 4.37 13.93 ± 0.74 11.31 ± 0.49	$\begin{array}{c} 1.86 \pm 0.08 \\ 10.31 \pm 0.19 \\ 5.24 \pm 0.41 \end{array}$	-2.13 ± 2.78 11.30 ± 1.06 9.64 ± 1.30	1.57 ± 1.78 13.50 ± 0.56 5.96 ± 0.56
$\gamma = 0.99$	HalfCheetah-v2 Hopper-v2 Walker2d-v2	511.99 ± 0.00 275.81 ± 0.00 346.63 ± 0.00	$\begin{array}{c} 137.30 \pm 70.70 \\ 155.19 \pm 16.27 \\ 244.67 \pm 30.28 \end{array}$	465.37 ± 4.56 276.10 ± 0.15 345.87 ± 1.79	-148.64 ± 10.50 238.63 ± 7.24 129.21 ± 16.38	-193.40 ± 88.66 227.96 ± 18.18 176.27 ± 22.33	251.49 ± 22.00 263.17 ± 3.47 338.60 ± 9.28	315.21 ± 57.95 178.17 ± 106.70 249.95 ± 27.41	305.69 ± 94.93 259.35 ± 5.48 314.11 ± 34.82
$\gamma = 0.999$	HalfCheetah-v2 Hopper-v2 Walker2d-v2	4097.30 ± 0.00 2223.49 ± 0.00 3151.77 ± 0.00	536.68 ± 384.66 408.25 ± 222.42 995.05 ± 330.00	3730.81 ± 31.57 1903.02 ± 83.18 2963.82 ± 26.59	-1150.45 ± 235.64 1878.19 ± 122.06 1039.71 ± 231.21	-1509.39 ± 877.41 1731.93 ± 233.34 1765.62 ± 212.79	3338.52 ± 191.06 2177.76 ± 43.64 2912.35 ± 335.22	2670.44 ± 437.00 1184.20 ± 805.75 1565.15 ± 1006.01	3303.42 ± 262.79 2187.72 ± 17.90 1877.53 ± 651.77

E.2 Imitating Environments

To evaluate algorithms for imitating environments, we add necessary information (e.g., robot position information) to the original state space defined by OpenAI Gym [10]. This is important since we need the learned environment model to predict the position information, upon which we can compute rewards for policy evaluation in the learned environments. Followed prior works [31, 25], the true reward function is assumed to be known in advance. We also normalize the robot position information by dividing 10 (but the reward function is not normalized). We use SAC [22] to re-train a sub-optimal policy as what we have done when imitating policies. We collect samples using this sub-optimality on true environments. Algorithmic configuration for BC and GAIL is the same as the one of imitating policies. Different from imitating-policies, the model output space (action space) is not bounded between -1 and +1. To overcome this difficulty, we normalize the model's outputs with statistics obtained from given demonstrations. During the training process, we also periodically evaluate the policy value of data-collecting policies on the learned environment models.

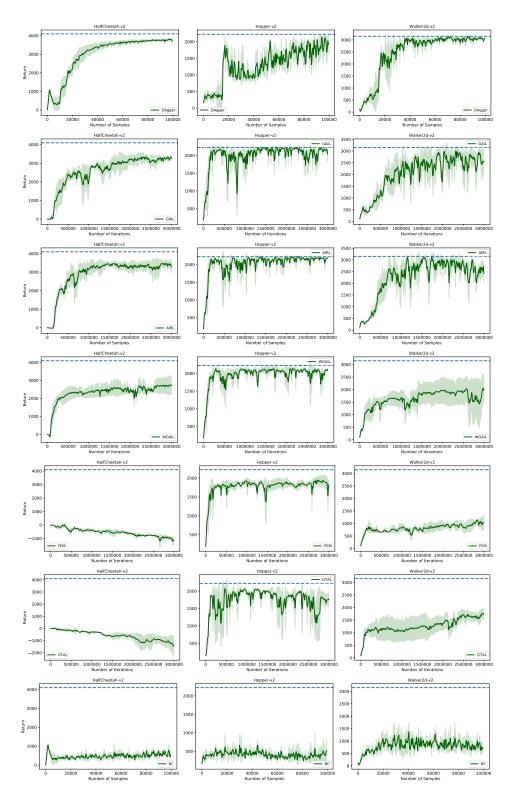


Figure 7: Learning curves of imitation approaches ($\gamma=0.999$) including DAgger, GAIL, AIRL, WGAIL, FEM, GTAL, and BC. The solid lines are mean of results and the shaded region corresponds to the stand deviation over 3 random seeds, while the dashed lines indicate the performance of expert policies.