

Distributionally Robust Behavioral Cloning for Robust Imitation Learning

Kishan Panaganti^{*,1}, Zaiyan Xu^{*,1}, Dileep Kalathil¹, Mohammad Ghavamzadeh²

Abstract—Robust reinforcement learning (RL) aims to learn a policy that can withstand uncertainties in model parameters, which often arise in practical RL applications due to modeling errors in simulators, variations in real-world system dynamics, and adversarial disturbances. This paper introduces the robust imitation learning (IL) problem in a Markov decision process (MDP) framework where an agent learns to mimic an expert demonstrator that can withstand uncertainties in model parameters without additional online environment interactions. The agent is only provided with a dataset of state-action pairs from the expert on a single (nominal) dynamics, without any information about the true rewards from the environment. Behavioral cloning (BC), a supervised learning method, is a powerful algorithm to address the vanilla IL problem. We propose an algorithm for the robust IL problem that utilizes distributionally robust optimization (DRO) with BC. We call the algorithm DR-BC and show its robust performance against parameter uncertainties both in theory and in practice. We also demonstrate the empirical performance of our approach to addressing model perturbations on several MuJoCo continuous control tasks.

Index Terms—Imitation Learning, Reinforcement Learning, Robust Reinforcement Learning

I. INTRODUCTION

A child, a dog, or even a reptile is capable of learning through imitation [1]. Such intuitive way of learning naturally extends from animal’s survival instincts to solving potentially complicated control tasks. Hence, it serves as the primary philosophy underlying most, if not all, methods in Imitation Learning (IL), a very fundamental reinforcement learning (RL) setting in which the goal is to learn a control policy exclusively from expert demonstrations. However simple and fundamental the idea of imitation learning may sound, variations in the training (simulators) and testing (real-world) environments can result in significant failures of current RL and IL control policies [2]–[5]. The training and testing environments in RL can vary due to several factors, such as modeling errors, real-world parameter changes, and adversarial disturbances. For instance, the sensor noise, action delay, friction, and mass of a mobile robot in the simulator may differ from those in the real-world setting. Furthermore, environmental conditions such as terrain, weather, lighting, and obstacle densities can also vary between the two settings which ultimately may make it infeasible to deploy learned policies to the real world.

Imitation Learning: Learning through imitation can be traced back to as early as [6]. Imitation learning assumes

access to only expert demonstrations. This has given rise to the most natural approach, behavior cloning [7], which is a supervised learning method and learns by simply minimizing differences between the actions of the learners and those of the experts for the states seen by the expert. [8] studied behavior cloning and characterized a tight bound on its sub-optimality gap of order $O(\varepsilon H^2)$. Most of works in imitation learning then try to improve this bound with additional assumptions. [9] proposed DAGGER which required hindsight expert actions for the states visited by the learner. They showed that the above bound on the sub-optimality gap can be improved to $O(\varepsilon uH)$, where u is the cost of taking a different action than the expert’s at one step and following the expert’s suggestions afterward, by querying the expert and interacting with the environment. In worse cases, u can be as large as H . Another approach, DRIL [10], needs environment interactions but no expert query, and achieves a sub-optimality gap bound that is linear in H . GAIL [11] uses a discriminator network to distinguish between the expert’s states from those visited by the learner’s policy while interacting with the environment. GAIL also achieves a sub-optimality gap bound that is linear in H [12]. In [13], the authors provide a game-theoretic framework for the IL problem that naturally competes with noisy expert policies. Different from these works, we focus on setting up an IL problem to address the parameter mismatch between the training and testing environments, and provide a practical algorithm to solve it.

Inverse RL (IRL): This is a framework where an agent learns the underlying true reward function of an expert and uses it in the usual RL algorithm to produce a policy that imitates the expert. Notable works include [14]–[18]. Although this is not the setting we consider, these are works that solve the IL problem using inverse RL methodologies. [19] proposed MIMIC-MD which directly estimates the expert trajectory distribution, but their algorithm is not practically implementable as is. [20] proposed MILO which uses an additional offline dataset to estimate environment dynamics. It has great empirical performance when trained with extremely limited expert demonstration, a scenario where BC fails to produce any working policy. [21] proposed an IRL algorithm to learn cost functions that are robust in noisy systems. These are all inherently adversarial approaches which are using critics to disturb the underlying systems or to distinguish and pick good reward representations. Different from these works, we focus on the min-max setting for IL where the objective is to learn the policy that minimizes the loss to mimic the expert’s actions against the worst possible models that lie in an uncertainty set.

^{*} Equal contributions, ¹ Authors are with the Department of Electrical and Computer Engineering at Texas A&M University, College Station, TX, USA. Email: {kpb, zxu43, dileep.kalathil}@tamu.edu, ² Author is with Google Research, USA. Email: ghavamza@google.com

Robust RL: The framework of the robust Markov decision process (RMDP) [22], [23] addresses the problem of learning a policy that is robust against mismatches between the training and testing environments. This is the goal of distributionally robust reinforcement learning (DR-RL). Robust RL has application in many real-world evolving systems in which there is always a gap between the true model and the simulator. Deploying naive RL policies [24] can be catastrophic when this gap is large. The RMDP problem is well-studied. [25]–[29] have investigated various types of uncertainty sets and sought tractable methods to solve RMDP. [30]–[32] have studied the sample complexity of model-based robust RL algorithms in a tabular setting using a generative model, which is a strong oracle enabling learners to query arbitrary transitions. [33] developed a model-free online robust RL algorithm with linear function approximation to tackle potentially infinite state spaces and [34] similarly developed a model-free *offline* robust RL algorithm with general function approximation. [35] proposed an online robust Q-learning with an R-contamination uncertainty set. We would like to note that robust RL has a strong connection to distributionally robust optimization (DRO). Many of the optimization techniques and analyses in robust RL were originally developed in the context of supervised learning by the DRO community [36]–[41]. This line of work is closest to ours, but as per our knowledge, ours is the first work to focus on addressing the parameter mismatch between the training and testing environments using DRO techniques in the context of *imitation learning* setting.

Main Contributions: We summarize our contributions in this paper as follows and refer to the relevant sections: **(i)** We introduce the problem of robust imitation learning for mismatch in model parameters. In this work, we consider mismatch in system transition dynamics. Robust learning in RL is studied widely but not in IL. We only know of [42], but it is in the IRL framework for making fair comparison. Critical real-world applications, such as power systems, healthcare, self-driving automobiles, have guidance from expert across diverse scenarios [43]–[45]. **(ii)** We propose a novel robust IL algorithm, called Distributionally Robust Behavioral Cloning (DR-BC). We provide theoretical guarantees for DR-BC. The BC method is computationally efficient which makes any other non-robust IL algorithm falls short on, and the DRO methodology addresses the model mismatch. Our proposed method cleverly bridges these two methods to solve the robust IL problem. We discuss this further in Section II. **(iii)** We perform extensive simulations on four notable continuous-action OpenAI [46] Gym MuJoCo [47] environments. We demonstrate that the DR-BC policy is robust against model perturbations in that when BC has catastrophic drop in performance, DR-BC weathers model mismatches for much more severe model perturbations.

II. ROBUST IMITATION LEARNING

In this section, we formally introduce our imitation learning problem that addresses parameter mismatch between the true and simulated models (transition dynamics). We then propose a robust IL algorithm for this problem that uses DRO, and

provide its theoretical guarantees. We end this section by discussing why we need to study robust IL.

A. Problem Formulation

The goal of IL in sequential decision-making is to *imitate* an expert’s policy by only using demonstrations generated by its interactions with the environment [10], [48]. More formally, consider an infinite-horizon MDP, denoted by the tuple $\{\mathcal{S}, \mathcal{A}, P^o, \gamma, r, \mu\}$, where \mathcal{S} and \mathcal{A} are the state and action spaces, $P^o : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition dynamics (*model*) of the environment, γ is a discount factor, μ is the initial state distribution, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the true reward function (unknown to the learner). In this paper, we consider a system with finite actions and a large state space. A stochastic policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ maps states to distributions over actions. For any policy π , the value function of an initial state $s_0 \sim \mu$ is given by $V_\pi(s_0) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | a_t \sim \pi(\cdot | s_t), s_{t+1} \sim P_{s_t, a_t}^o]$. We denote $V_\pi = \mathbb{E}_{s_0 \sim \mu} V_\pi(s_0)$ and drop the explicit dependence on μ going forward for notation simplicity.

For any policy π , denote $d_\pi^o \in \Delta(\mathcal{S})$ as the state distribution of π under the evaluation of model P^o with initial state picked from μ . In this section, we simply denote such state-distributions as d^π and make explicit dependence on P^o where brevity is needed. Formally, let $\Pr_t(s | \pi, s_0 \sim \mu)$ be the probability of visiting state $s \in \mathcal{S}$ at time t following policy π on model P^o starting at initial state $s_0 \sim \mu$. Then, the state distribution of π is $d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr_t(s | \pi, s_0 \sim \mu, P^o)$. We can now rewrite the value of π as $V_\pi = \mathbb{E}_{s \sim d^\pi, a \sim \pi} [r(s, a)] / (1 - \gamma)$.

In the vanilla IL setting, the true reward function r is unknown to the learner. We instead have the dataset generated by rolling an expert policy (which is unknown to the learner) specified by $\pi_e : \mathcal{S} \rightarrow \Delta(\mathcal{A})$. Concisely speaking, we have an expert dataset in the form of i.i.d. tuples $\mathcal{D}_e = \{s_i, a_i\}_{i=1}^N$ sampled from state distribution $d_{P^o}^{\pi_e}$ and an expert policy π_e .

We use the RMDP framework [22], [23] subsuming the MDP framework described above. Consider an RMDP tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{P}, \gamma, r, \mu\}$ where $\gamma \in [0.5, 1)$ and the *uncertainty set* \mathcal{P} is defined as

$$\mathcal{P} = \otimes_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_{s,a} \quad \text{with} \\ \mathcal{P}_{s,a} = \{P_{s,a} \in \Delta(\mathcal{S}) : D(P_{s,a}, P_{s,a}^o) \leq \rho'_r\}, \quad (1)$$

where $P^o = (P_{s,a}^o, (s, a) \in \mathcal{S} \times \mathcal{A})$ is the simulator model, $D(\cdot, \cdot)$ is a distance measure between two probability distributions (e.g., total variation, chi-square, Kullback-Liebler (f -divergences in general), Wasserstein), and $\rho'_r \in (0, (1 - \gamma)/\gamma]$ is the radius of the uncertainty set indicating the level of robustness. We assume the real-world model belongs to this uncertainty set \mathcal{P} . We restrict to the total-variation distance D_{TV} for the measure D in this paper and leave other types of measures for future work.

From the RMDP literature [22], [32], [34], [49], we introduce the *robust* value function as $V_\pi^{\text{rob}}(s) = \sum_a \pi(a | s) Q_\pi^{\text{rob}}(s, a)$ and the corresponding *robust* Q-value function as $Q_\pi^{\text{rob}}(s, a) = r(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{P}_{s,a}} P_{s,a}^\top V_\pi^{\text{rob}}$ for policy π . Similar to [8], we let $\pi_e \in \Pi$, where Π is the class of

stochastic policies, be a *good robust* policy (under the above RMDP setting). That is, it satisfies $\max_{\pi \in \Pi} V_{\pi}^{\text{rob}} - V_{\pi_e}^{\text{rob}} \leq o(H)$ (something small compared to time horizon). For notation simplicity, as in [8], [20], [48], we just let $V_{\pi_e}^{\text{rob}} \geq V_{\pi}^{\text{rob}}$ hold for all $\pi \in \Pi$. **Now, we pose the robust IL problem as follows.** The goal of a robust IL algorithm is to output a policy $\hat{\pi}$ that imitates the expert policy π_e by satisfying $V_{\hat{\pi}}^{\text{rob}} \approx V_{\pi_e}^{\text{rob}}$. We have provided real-world applications that motivate this problem formulation in Section I.

B. Need for Robust Imitation Learning

In this section, we formally show that the vanilla behavioral cloning policy Equation (2) can be arbitrarily bad (as bad as a random policy) compared to an expert policy (*good robust* policy). Consider the following vanilla behavioral cloning [8] optimization problem.

$$\pi_{\text{bc}} = \arg \min_{\pi} L_{\text{bc}}(\pi) = \mathbb{E}_{s \sim d^{\pi_e}} [l(\pi_e(\cdot | s), \pi(\cdot | s))]. \quad (2)$$

We assume access to sampling possibly *infinite data* from the state distribution d^{π_e} to calculate the loss $L_{\text{bc}}(\pi_{\text{bc}})$ up to some small error. We consider a simple setting with $\mathcal{P} = \{P^o, P'\}$ where P^o is the simulator model and P' is the perturbed model. We give the following result similar to [50, Theorem 4] and skip its proof.

Theorem 1 (Robustness Gap). *There exists an uncertainty set $\mathcal{P} = \{P^o, P'\}$, initial state $s_0 \in \mathcal{S}$, expert policy π_e such that $\max_{\pi \in \Pi} V_{\pi}^{\text{rob}}(s_0) - V_{\pi_e}^{\text{rob}}(s_0) \leq \varepsilon$ for small $\varepsilon > 0$, and discount factor $\gamma \in (\gamma_o, 1]$ such that $V_{\pi_{\text{bc}}}^{\text{rob}}(s_0) \leq V_{\pi_e}^{\text{rob}}(s_0) - c/(1 - \gamma)$, where c is a positive constant.*

Remark 1. The vanilla behavioral cloning policy π_{bc} compared to expert policy (*good robust* policy) is bad with a performance gap $\Omega(1/(1 - \gamma))$. Since $|r(s, a)| \leq 1$ uniformly by assumption, $\|V_{\pi}^{\text{rob}}\|_{\infty} \leq 1/(1 - \gamma)$ for any policy π . Therefore, the difference between the optimal/expert robust value function and the robust value function of an arbitrary policy cannot be greater than $\mathcal{O}(1/(1 - \gamma))$. Thus the performance of π_{bc} can be as bad as an arbitrary policy in an order sense. In the next section, we propose an algorithm to solve the robust imitation learning problem.

C. Robust Against Model Mismatch

We propose a principled adversarial approach by the methodology of distributionally robust optimization (DRO) to solve the robust imitation learning problem. DRO is now a well-established area [36], [40], [51], whose formulation is identical to that in the classical RMDP [22], [23] in DR-RL. The distributionally robust behavioral cloning algorithm solves the following optimization problem getting the policy π_{drbc} :

$$\arg \min_{\pi} \max_{P \in \mathcal{P} : D_{\text{TV}}(d_P^{\pi_e}, d_{P^o}^{\pi_e}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} [l(\pi(\cdot | s), \pi_e(\cdot | s))], \quad (3)$$

where ρ_r , the *robustness radius parameter* which is a problem-dependent constant, is set to $\gamma \rho_r' / (1 - \gamma) \in (0, 1]$, and $l(\pi(\cdot | s), \pi_e(\cdot | s))$ is a surrogate loss function which measures how far the learner policy π is with respect to the expert action for the states visited by the expert. Examples

of the loss function l comprise of 0-1 loss (described by $\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbf{1}(a \neq \pi_e(s))$ for deterministic expert policies), total variation loss (described by $D_{\text{TV}}(\pi_e(\cdot | s), \pi(\cdot | s)) = 0.5 \|\pi_e(\cdot | s) - \pi(\cdot | s)\|_1$), KL loss (described by $D_{\text{KL}}(\pi(\cdot | s) \|\pi_e(\cdot | s)) = \sum_a \pi(a|s) \log(\pi(a|s)/\pi_e(a|s))$ with π absolutely continuous to π_e), and many more such quantifiers. We simply use the D_{TV} loss function in this paper considering its known connections with other f -divergences [52], [53]. We note that the DR-BC policy π_{drbc} depends on ρ_r but we simply choose to make it inherent for notation simplicity. We also remark that we recover the behavioral cloning policy [7] with $\rho_r = 0$ in DR-BC policy Equation (3).

We define the uncertainty set parameterized by ρ_r as $\mathcal{M} = \{P \in \mathcal{P} : D_{\text{TV}}(d_P^{\pi_e}, d_{P^o}^{\pi_e}) \leq \rho_r\}$. It is straightforward from its definition and Lemma 7 that $\mathcal{M} = \mathcal{P}$. The DR-BC algorithm Equation (3) finds π_{drbc} for the IL problem by minimizing an observed surrogate loss between its actions and the actions of an expert policy under the adversarial state distribution for a model in class \mathcal{M} which acts as a worse-case distribution. We define the model mismatch distributionally robust behavioral cloning loss function as $L_{\text{drbc}}(\pi, \rho_r) = \max_{P \in \mathcal{M}} \mathbb{E}_{s \sim d_P^{\pi_e}} [D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s))]$ for any policy π and ρ_r . But we immediately notice that to solve the inner optimization in Equation (3) we need access to all the state distributions around the expert's state distribution. Even knowing the model P^o , this is computationally intractable. Moreover, we would also need the capability of querying an expert for actions for various states chosen by such state distributions. Also assuming having access to all models in \mathcal{M} is unrealistic. We now discuss circumventing this challenge using the DRO methodology [36], [37].

Motivated from the DR-RL literature [32], [34], [49], we now have the following result that provides a dual reformulation for the inner maximization in Equation (3) as a consequence of the DRO methodology.

Proposition 2. *For a fixed expert policy $\pi_e \in \Pi$, we have, for all $\pi \in \Pi$ and $\rho_r \in (0, 1]$,*

$$\begin{aligned} & \max_{P \in \mathcal{M}} \mathbb{E}_{s \sim d_P^{\pi_e}} [D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s))] \\ &= \min_{\eta \in \mathbb{R}} \mathbb{E}_{s \sim d_{P^o}^{\pi_e}} [(D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s)) - \eta)_+] \\ & \quad + \left(\sup_{s \in \mathcal{S} : d_{P^o}^{\pi_e}(s) > 0} D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s)) - \eta \right)_+ \cdot \rho_r + \eta. \end{aligned}$$

Proof. We first rewrite $\max_{P \in \mathcal{M}} \mathbb{E}_{s \sim d_P^{\pi_e}} [D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s))]$ as $\max_{d_P^{\pi_e} : D_{\text{TV}}(d_P^{\pi_e}, d_{P^o}^{\pi_e}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} [D_{\text{TV}}(\pi(\cdot | s), \pi_e(\cdot | s))]$ since $\mathcal{M} = \{P \in \mathcal{P} : D_{\text{TV}}(d_P^{\pi_e}, d_{P^o}^{\pi_e}) \leq \rho_r\}$. Then the statement immediately follows from Lemma 5. \square

We give our DR-BC algorithm that only requires an expert dataset \mathcal{D}_e generated according to model P^o in Algorithm 1 based on Proposition 2. The DRO technique in Proposition 2 transforms the inner maximization in Equation (3) to an unconstrained scalar variables convex optimization problem. We remark that this new optimization problem due to the dual reformulation only depends on the expert's state distribution. This enables us to use the expert dataset to solve the DR-BC objective Equation (3). We emphasize that we need access to

all the state distributions to solve the inner optimization in Equation (3) directly which is computationally intractable for large-scale problems. Now we are overcoming this challenge through this dual reformulation result. We refer to Section IV for further details.

Algorithm 1 Distributionally Robust Behavioral Cloning

- 1: **Input:** Expert dataset $\mathcal{D}_e = (s_i, a_i)_{i=1}^N$ according to model P^o , model mismatch radius parameter ρ_r .
- 2: **Initialize:** Policy π_θ parameterized by θ .
- 3: Calculate the empirical loss for $L_{\text{drbc}}(\pi_\theta, \rho_c)$:

$$\min_{\eta \in \mathbb{R}} \left(\frac{1}{N} \sum_{(s,a) \in \mathcal{D}_e} (l(a, \pi_\theta(s)) - \eta)_+ \right) + \rho_r \left(\sup_{(s,a) \in \mathcal{D}_e} l(a, \pi_\theta(s)) - \eta \right)_+ + \eta. \quad (4)$$

- 4: $\theta \leftarrow \arg \min_{\theta} L_{\text{drbc}}(\pi_\theta, \rho_r)$.
 - 5: **Output policy:** $\hat{\pi}_{\text{drbc}} = \pi_\theta$
-

We now give the sub-optimality guarantee of model mismatch DR-BC policy. We provide its proof in Section III.

Theorem 3 (Model mismatch DR-BC sub-optimality bound). *Assume small optimization error $L_{\text{drbc}}(\pi_{\text{drbc}}, \rho_r) = \varepsilon_{\text{drbc}}(\rho_r)$. We have $V_{\pi_e}^{\text{rob}} - V_{\pi_{\text{drbc}}}^{\text{rob}} \leq 2\varepsilon_{\text{drbc}}(\rho_r)/(1-\gamma)^2$.*

Remark 2. We have an $\mathcal{O}(\varepsilon_{\text{drbc}}(\rho_r)H^2)$ sub-optimality bound from Theorem 3. With a small optimization error $\varepsilon_{\text{drbc}}(\rho_r)$, the sub-optimality guarantee for DR-BC algorithm is superior to the BC policy as discussed in Section II-B. When the robustness parameter $\rho_r = 0$, we recover the non-robust BC algorithm and its quadratic horizon dependence [8]. This sub-optimality bound is in fact tight $\Omega(\varepsilon H^2)$ [8], [48].

We also present the approximation result for the sub-optimality of $\hat{\pi}_{\text{drbc}}$ returned by Algorithm 1 that uses the expert dataset \mathcal{D}_e . We provide its proof in Section III. We consider e_{\min} , the minimum non-zero probability value in π_e , as a problem dependent constant. We again consider $L_{\text{drbc}}(\hat{\pi}_{\text{drbc}}, \sqrt{\rho_r}) = \varepsilon_{\text{drbc}}(\rho_r) > 0$ be a small optimization error for all $\rho_r \in (0, 1]$.

Theorem 4 (Approximate DR-BC sub-optimality bound). *Let $\hat{\varepsilon}_{\text{drbc}}(\rho_r) = \varepsilon_{\text{drbc}}(\rho_r) + \tilde{\mathcal{O}}(\rho_r \sqrt{\log(1/\delta)} / (e_{\min} N))$. Then, for any $\rho_r \in (0, 1]$, policy $\hat{\pi}_{\text{drbc}}$ satisfies $V_{\pi_e}^{\text{rob}} - V_{\hat{\pi}_{\text{drbc}}}^{\text{rob}} \leq 2\hat{\varepsilon}_{\text{drbc}}(\rho_r)/(1-\gamma)^2$, with probability at least $1 - \delta$.*

Remark 3. We note that $\tilde{\mathcal{O}}(\cdot)$ is order optimal up to a logarithmic term on N and its exact form is available in Section III. The approximate sub-optimality guarantee for DR-BC algorithm is still superior to the BC policy as discussed in Section II-B. We indeed showcase empirically as well that DR-BC algorithm is resilient to model perturbations in Section IV.

III. ANALYSIS DETAILS

A. Useful Technical Results

We now state a result from [34] based on DRO methodology which is useful for proving Proposition 2.

Lemma 5 ([34, Lemma 5]). *Let P^o be a non-zero distribution on the space \mathcal{X} and $l : \mathcal{X} \rightarrow \mathbb{R}$ be a loss function. Then*

$$\sup_{D_{\text{TV}}(P, P^o) \leq \rho} \mathbb{E}_{x \sim P}[l(x)] = \inf_{\eta \in \mathbb{R}} \left\{ \mathbb{E}_{x \sim P^o}[(l(x) - \eta)_+] + (\sup_{x \in \mathcal{X}} l(x) - \eta)_+ \cdot \rho + \eta \right\}. \quad (5)$$

We now specialize [36, Corollary 2] for the total variation distance.

Lemma 6. *Let $\Theta \subseteq \mathbb{R}^d$, $l : \mathcal{X} \times \Theta \mapsto [0, M]$ and fix any $\rho \in (0, 1]$. We have*

$$\sup_{D_{\text{TV}}(P, P^o) \leq \rho} \mathbb{E}_P[l(X, \hat{\theta})] \leq \inf_{\theta \in \Theta} \sup_{D_{\text{TV}}(P, P^o) \leq \sqrt{\rho}} \mathbb{E}_P[l(X, \theta)] + cM(1 + \rho)\sqrt{(\log(1/\delta) + 2d \log(N))/N},$$

which holds with probability at least $1 - \delta$, where $c > 0$ is some universal constant and $\hat{\theta} = \arg \min_{\theta \in \Theta} \sup_{D_{\text{TV}}(P, P^o) \leq \sqrt{\rho}} \mathbb{E}_P[l(X, \theta)]$.

Proof. The proof simply follows from observation $\{p : D_{\text{TV}}(p, q) \leq 2\rho'\} \subseteq \{p : D_{\chi^2}(p, q) \leq \rho'\} \subseteq \{p : D_{\text{TV}}(p, q) \leq 2\sqrt{\rho'}\}$ which follows from Pinsker's inequality [52], [53, Theorem 5], and [54, Lemma 11.1]. \square

B. Proof of Theorem 3

We present a few results needed for proving Theorem 3. First, we formally show that when two models are close, then their state-distributions are close under the same policy.

Lemma 7. *Consider any policy π and $P \in \mathcal{P}$. We have $D_{\text{TV}}(d_P^\pi, d_{P^o}^\pi) \leq \gamma\rho_r/(1-\gamma)$.*

Proof. By definition, since $P \in \mathcal{P}$, we have $D_{\text{TV}}(P_{s,a}, P_{s,a}^o) \leq \rho_r'$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. We denote matrices $\mathbb{P}_\pi, \mathbb{P}_\pi^o : \mathcal{S} \times \mathcal{S} \rightarrow [0, 1]$ with $\mathbb{P}_\pi(s', s) = \sum_{a \in \mathcal{A}} \pi(a|s)P_{s,a}(s')$ and $\mathbb{P}_\pi^o(s', s) = \sum_{a \in \mathcal{A}} \pi(a|s)P_{s,a}^o(s')$. Now, we can write

$$\begin{aligned} d_P^\pi &= (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr_t(s \mid \pi, s_0 \sim \mu, P) \\ &= (1-\gamma) \sum_{t=0}^{\infty} (\gamma \mathbb{P}_\pi)^t \mu, \end{aligned}$$

and similarly $d_{P^o}^\pi = (1-\gamma) \sum_{t=0}^{\infty} (\gamma \mathbb{P}_\pi^o)^t \mu$. Denoting $\mathbb{P}_{t,\pi} = \mathbb{P}_\pi^t \mu, \mathbb{P}_{t,\pi}^o = (\mathbb{P}_\pi^o)^t \mu$, from triangle inequality we further get

$$\|d_P^\pi - d_{P^o}^\pi\|_1 \leq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \|\mathbb{P}_{t,\pi} - \mathbb{P}_{t,\pi}^o\|_1. \quad (6)$$

Intuitively, $\mathbb{P}_{t,\pi}(\mathbb{P}_{t,\pi}^o)$ is state distribution resulting from π evolving in the model $P(P^o)$ at time step t with μ as the initial state distribution. We now bound $\|\mathbb{P}_{t,\pi} - \mathbb{P}_{t,\pi}^o\|_1$ for $t \geq 0$ in a recursive approach. From basic Markov chain

theory [55] for any $t \geq 0$, we have

$$\begin{aligned}
& \|\mathbb{P}_{t,\pi} - \mathbb{P}_{t,\pi}^o\|_1 = \sum_{s'} |\mathbb{P}_{t,\pi}(s') - \mathbb{P}_{t,\pi}^o(s')| \\
&= \sum_{s'} \left| \sum_{s,a} (\mathbb{P}_{t-1,\pi}(s) P_{s,a}(s') - \mathbb{P}_{t-1,\pi}^o(s) P_{s,a}^o(s')) \pi(a|s) \right| \\
&\leq \sum_s |\mathbb{P}_{t-1,\pi}(s) - \mathbb{P}_{t-1,\pi}^o(s)| \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) \\
&\quad + \sum_s \mathbb{P}_{t-1,\pi}^o(s) \sum_a \pi(a|s) \sum_{s'} |P(s'|s,a) - P^o(s'|s,a)| \\
&\leq \|\mathbb{P}_{t-1,\pi} - \mathbb{P}_{t-1,\pi}^o\|_1 + 2\rho_r',
\end{aligned}$$

where the last inequality holds since $D_{\text{TV}}(P_{s,a}, P_{s,a}^o) = (1/2) \|P_{s,a} - P_{s,a}^o\|_1 \leq \rho_r'$ for all $(s,a) \in \mathcal{S} \times \mathcal{A}$. By recursion, we have $\|\mathbb{P}_{t,\pi} - \mathbb{P}_{t,\pi}^o\|_1 \leq 2\rho_r' t$. Recall from algebra that $\sum_{t=0}^{\infty} \gamma^t t = \gamma/(1-\gamma)^2$. Combining this with Equation (6) completes the proof. \square

Now we state a result which extends the performance difference lemma [56, Lemma 1.16] notion for robust MDPs.

Lemma 8 (Robust Performance Difference Lemma). *For any π', π policies, we get*

$$\begin{aligned}
V_{\pi}^{\text{rob}} - V_{\pi'}^{\text{rob}} &\leq \frac{1}{1-\gamma} \\
&\max_{P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} \left[\sum_a (\pi(a|s) - \pi'(a|s)) Q_{\pi'}^{\text{rob}}(s,a) \right].
\end{aligned}$$

Proof. We first define few useful notations for this proof. The *robust model* $P^{\text{rob},\pi}$ for every π is as follows: $P_{s,a}^{\text{rob},\pi} = \arg \min_{P_{s,a} \in \mathcal{P}_{s,a}} P_{s,a}^{\top} V_{\pi}^{\text{rob}}$, $(s,a) \in \mathcal{S} \times \mathcal{A}$. We call V_{π}^P as the value function for policy π under the model P . Now we can write $V_{\pi}^{\text{rob}} = V_{\pi'}^{P^{\text{rob},\pi}}$.

Fix $s_0 \sim \mu$. For any π', π policies, we have

$$\begin{aligned}
V_{\pi}^{\text{rob}}(s_0) - V_{\pi'}^{\text{rob}}(s_0) &\stackrel{(a)}{\leq} V_{\pi}^{P^{\text{rob},\pi'}}(s_0) - V_{\pi'}^{\text{rob}}(s_0) \\
&\stackrel{(b)}{=} V_{\pi}^{P^{\text{rob},\pi'}}(s_0) - V_{\pi'}^{P^{\text{rob},\pi'}}(s_0) \\
&\stackrel{(c)}{=} \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_e}, P^{\text{rob},\pi'}} \left[\sum_a (\pi(a|s) - \pi'(a|s)) Q_{\pi'}^{\text{rob}}(s,a) \right],
\end{aligned}$$

where (a) follows since by definition of $V_{\pi}^{\text{rob}}(s_0)$ we have $V_{\pi}^{\text{rob}}(s_0) \leq V_{\pi}^{P^{\text{rob},\pi'}}(s_0)$, and (b) follows from definition of $P^{\text{rob},\pi'}$ yielding $V_{\pi'}^{\text{rob}}(s_0) = V_{\pi'}^{P^{\text{rob},\pi'}}(s_0)$. Observe $P^{\text{rob},\pi'} \in \mathcal{P}$. Now taking expectation on $s_0 \sim \mu$ with Lemma 7 completes the proof of this result. Now it only remains to show (c).

For (c), first denote $\mathcal{T}_{\pi,\pi'}(s) = (s_t, a_t)_{t \geq 0}$ trajectory generated from rolling policy π from the initial state s_0 under the robust model $P^{\text{rob},\pi'}$. Now,

$$\begin{aligned}
& V_{\pi}^{P^{\text{rob},\pi'}}(s_0) - V_{\pi'}^{P^{\text{rob},\pi'}}(s_0) \\
&= \mathbb{E}_{\mathcal{T}_{\pi,\pi'}(s)} \left[\sum_t \gamma^t r(s_t, a_t) \right] - V_{\pi'}^{P^{\text{rob},\pi'}}(s_0)
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(d)}{=} \mathbb{E}_{\mathcal{T}_{\pi,\pi'}(s)} \left[\sum_t \gamma^t (r(s_t, a_t) + \gamma V_{\pi'}^{P^{\text{rob},\pi'}}(s_{t+1}) \right. \\
&\quad \left. - V_{\pi'}^{P^{\text{rob},\pi'}}(s_t) \right] \\
&\stackrel{(e)}{=} \mathbb{E}_{\mathcal{T}_{\pi,\pi'}(s)} \left[\sum_t \gamma^t (Q_{\pi'}^{\text{rob}}(s_t, a_t) - V_{\pi'}^{\text{rob}}(s_t)) \right] \\
&\stackrel{(f)}{=} \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{s_0}^{\pi_e}, P^{\text{rob},\pi'}} \sum_{a'} \pi(a'|s') (Q_{\pi'}^{\text{rob}}(s', a')) \\
&\quad - \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{s_0}^{\pi_e}, P^{\text{rob},\pi'}} \sum_{a'} \pi'(a'|s') (Q_{\pi'}^{\text{rob}}(s', a')),
\end{aligned}$$

where (d) follows by recursion, (e) follows since $Q_{\pi'}^{\text{rob}}(s,a) = r(s,a) + \gamma (P_{s,a}^{\text{rob},\pi'})^{\top} V_{\pi'}^{P^{\text{rob},\pi'}}$, and (f) from $V_{\pi'}^{\text{rob}}(s) = \sum_a \pi'(a|s) Q_{\pi'}^{\text{rob}}(s,a)$. This proves (c). \square

Proof of Theorem 3. We start by Lemma 8 with $\pi' = \pi_{\text{drbc}}$ and $\pi = \pi_e$. We get

$$\begin{aligned}
V_{\pi_e}^{\text{rob}} - V_{\pi'}^{\text{rob}} &\leq \frac{1}{1-\gamma} \\
&P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r \max_{P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} \left[\sum_a (\pi_e(a|s) - \pi'(a|s)) Q_{\pi'}^{\text{rob}}(s,a) \right] \\
&\stackrel{(a)}{\leq} \frac{1}{1-\gamma} \\
&P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r \max_{P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} \left[\|\pi_e(\cdot|s) - \pi'(\cdot|s)\|_1 \|Q_{\pi'}^{\text{rob}}(s, \cdot)\|_{\infty} \right] \\
&\stackrel{(b)}{\leq} \frac{2}{(1-\gamma)^2} P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r \max_{P: D(d_{P^e}^{\pi_e}, d_{P^o}^{\pi_o}) \leq \rho_r} \mathbb{E}_{s \sim d_P^{\pi_e}} [D_{\text{TV}}(\pi'(\cdot|s) \|\pi_e(\cdot|s))],
\end{aligned}$$

where (a) follows from Holder's inequality and (b) from D_{TV} definition and the fact that $\|Q_{\pi'}^{\text{rob}}(s, \cdot)\|_{\infty} \leq 1/(1-\gamma)$ for any π' . The proof of this result is complete since $L_{\text{drbc}}(\pi_{\text{drbc}}, \rho_r) = \varepsilon_{\text{drbc}}(\rho_r)$. \square

C. Proof of Theorem 4

Proof. Firstly, from Lemma 6, we observe that $L_{\text{drbc}}(\hat{\pi}_{\text{drbc}}, \rho_c) \leq \hat{\varepsilon}_{\text{drbc}}(\rho_c)$ holds with probability at least $1 - \delta$ with

$$\hat{\varepsilon}_{\text{drbc}}(\rho_c) = \varepsilon_{\text{drbc}}(\rho_c) + c(1 + \rho_c) \sqrt{\frac{\log(1/\delta) + 2|\mathcal{A}| \log(N)}{e_{\min} N}},$$

where $e_{\min} = \min_{s,a: \pi_e(a|s) > 0} \pi_e(a|s)$ and $c > 0$ is some universal constant. The proof is now complete by following the analysis of Theorem 3 with $\hat{\pi}_{\text{drbc}}$. \square

IV. EXPERIMENTS

We aim to answer the question: When model mismatches are present, is the DR-BC algorithm robust compared to the non-robust BC algorithm?

A. Experiment Setup and Practical Algorithm

We perform extensive simulations on four OpenAI Gym [46] environments simulated with MuJoCo physics engine [47]: Hopper-v3, HalfCheetah-v3, Walker2d-v3, and Ant-v3. We train both the BC and DR-BC algorithms on the expert data generated by the pre-trained TD3 [57] policies from the RL Baselines3 Zoo repositories [58]. [20] pointed out that BC is very effective at imitating the expert when

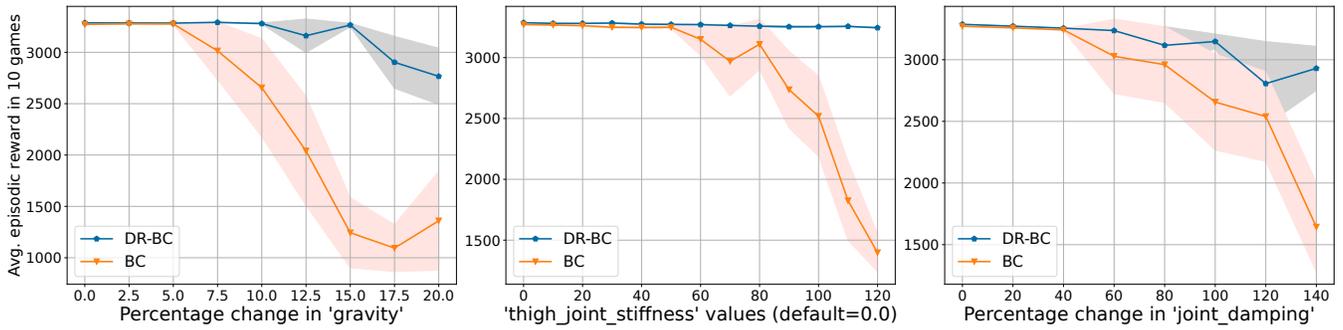


Fig. 1: Hopper-v3 *perturbation results*. Average episodic reward on 10 differently seeded episodes. From left to right, the perturbations are in: 'gravity', 'joint_stiffness' of the thigh joint, and 'joint_damping' of all joints.

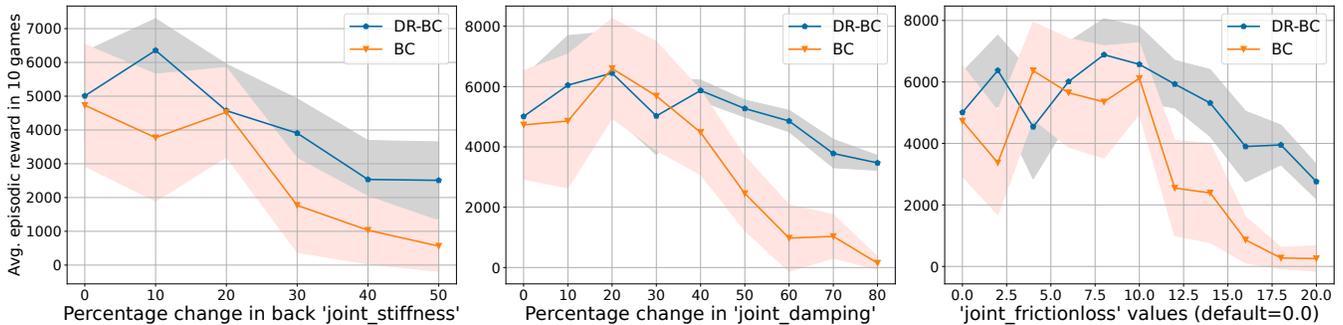


Fig. 2: HalfCheetah-v3 *perturbation results*. Average episodic reward on 10 differently seeded episodes. From left to right, the perturbations are in: 'joint_stiffness' of all back joints, 'joint_damping' of all joints, and 'joint_frictionloss' of all joints.

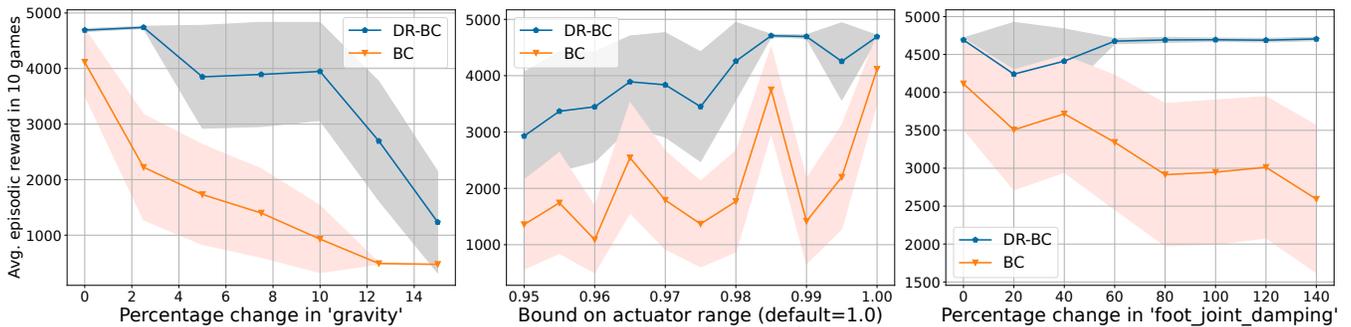


Fig. 3: Walker2d-v3 *perturbation results*. Average episodic reward on 10 differently seeded episodes. From left to right, the perturbations are in: 'gravity', 'actuator_ctrlrange' of all joints, and 'joint_damping' of both foot joints.

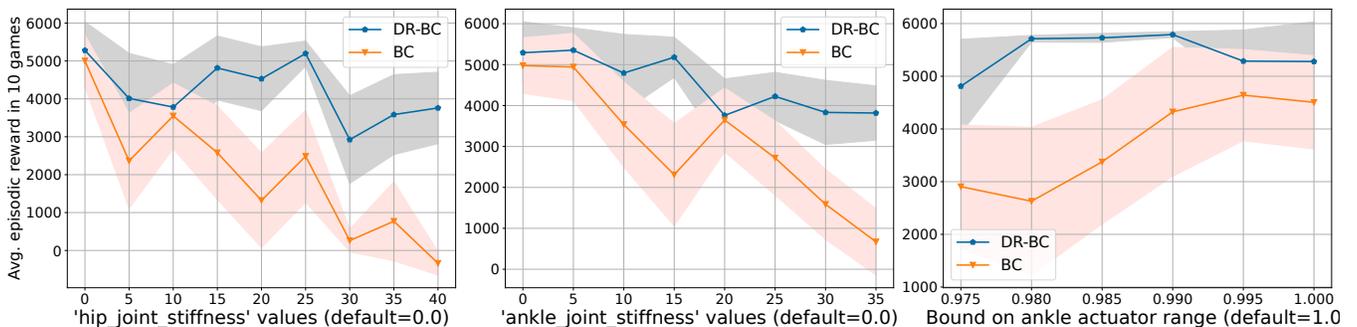


Fig. 4: Ant-v3 *perturbation results*. Average episodic reward on 10 differently seeded episodes. From left to right, the perturbations are in: 'joint_stiffness' of hip joints, 'joint_stiffness' of ankle joints, and 'actuator_ctrlrange' of ankle joints.

given large number of samples. Hence, like [10], [20], we give both BC and DR-BC the same but relatively low number of expert trajectories. *For reproducibility, we provide the code and more implementation details in the GitHub repository <https://github.com/ferocious-cheetah/DRBC>.*

Now we explain Algorithm 1 in details. We have the following proposition, similar to [34, Proposition 1] and skip its proof, showing that the unconstrained optimization of the dual variable in Equation (4) of Algorithm 1 can be further reduced to be over just a finite real interval.

Proposition 9. *Suppose that we have deterministic policies π_e , π and a bounded action space \mathcal{A} . Let the loss l be chosen as the squared L2 loss, i.e., $l(\pi(s), \pi_e(s)) = \|\pi(s) - \pi_e(s)\|_2^2$. Further, denote $\bar{L} = \sup_{(s,a) \in \mathcal{D}_e} \|\pi(s) - \pi_e(s)\|_2^2$. Then the dual reformulation in Equation (4) can be further rewritten as*

$$L_{\text{drbc}}(\pi_\theta, \rho_r) = \inf_{\eta \in [0, (1+\rho_r)\bar{L}]} \left\{ \frac{1}{N} \sum_{(s,a) \in \mathcal{D}_e} (l(a, \pi_\theta(s)) - \eta)_+ + \rho_r(\bar{L} - \eta)_+ + \eta \right\}. \quad (7)$$

At initialization, we need an expert dataset \mathcal{D}_e of size N and some radius of our uncertainty set ρ_r . We also need to initialize a neural network π_θ which is our policy (actor) with random parameters θ . In each iteration, to solve Equation (7), we use the minimization solver from the powerful optimization libraries in SciPy [59]. In particular, we use the SLSQP method [60] with the bounds prescribed in Proposition 9. In step 4, our policy (actor) is optimized based on the L_{drbc} loss using ADAM [61].

B. Test For Robustness

When the testing environment is perturbed, e.g., change in gravity, perturbed actuator and modified damping coefficient, model mismatches are present. Here we explain the simulation results on each of the four environments in details.

We perturb Hopper-v3 by changing the model parameter ‘gravity’, ‘thigh_joint_stiffness’, and ‘joint_damping’. Figure 1 shows that DR-BC is tenacious under model perturbations. For example, in the middle figure, when the ‘thigh_joint_stiffness’ parameter is positive and increasing, a joint spring is created in the thigh of hopper and becomes stiffer. A non-robust policy such as BC cannot withstand such mismatch between the training and testing environments. Meanwhile, our DR-BC agent refuses to drop in performance. In Figure 2 and Figure 3, DR-BC still refuses to lose performance in wide ranges of perturbations on different environments.

Ant-v3 is the most difficult environment among the four. In Figure 4, we perturb it by changing the model parameter ‘joint_stiffness’ of all four hip joints, ‘joint_stiffness’ of all four ankle joints, and ‘actuator_ctrlrange’ of all four ankle joints. When the control range of a joint actuator is reduced, it becomes harder for the agent to recover from dramatic change in the posture that involves that joint, let alone simultaneous perturbation in all four joint actuators. The rightmost figure in Figure 4 shows that the performance of BC precipitates

when it no longer has the full control range. On the other hand, DR-BC performance is stable throughout.

V. CONCLUSION

In this paper, we introduce a novel problem of *robust imitation learning* to incorporate resiliency to changes in the real-world parameters. We present a novel approach to solve this problem, Distributionally Robust Behavioral Cloning (DR-BC) algorithm. Our proposed DR-BC algorithm utilizes the distributionally robust optimization (DRO) technique for BC to efficiently address robustness for the changes in real-world parameters. We have shown through both theoretical and practical analysis that DR-BC can effectively and computationally efficiently combat the model perturbations in many benchmark MuJoCo tasks.

While in this paper we only consider the total variation distance for the inner maximization, future work will explore using other types measures such as KL-divergence and Chi-square divergence. The same applies to the loss function considered in this work. We also plan to work on the scenario where the model is not known in large-scale problems using general function approximations. An interesting practical direction could be to use DR-BC algorithm to fine-tune the policy network in online IL algorithms like GAIL which generate more diverse and realistic examples.

REFERENCES

- [1] A. Kis, L. Huber, and A. Wilkinson, “Social learning by imitation in a reptile (pogona vitticeps),” *Animal cognition*, vol. 18, 09 2014.
- [2] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Ucroft, P. Abbeel, W. Burgard, M. Milford, *et al.*, “The limits and potentials of deep learning for robotics,” *The International journal of robotics research*, vol. 37, no. 4-5, pp. 405–420, 2018.
- [3] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, “Domain randomization for transferring deep neural networks from simulation to the real world,” in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 23–30, 2017.
- [4] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-real transfer of robotic control with dynamics randomization,” in *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3803–3810, IEEE, 2018.
- [5] D. C. Guastella and G. Muscato, “Learning-based methods of perception and navigation for ground vehicles in unstructured environments: A review,” *Sensors*, vol. 21, no. 1, p. 73, 2020.
- [6] D. A. Pomerleau, “Alvinn: An autonomous land vehicle in a neural network,” in *Advances in Neural Information Processing Systems* (D. Touretzky, ed.), vol. 1, Morgan-Kaufmann, 1988.
- [7] M. Bain and C. Sammut, “A framework for behavioural cloning,” in *Machine Intelligence 15*, 1995.
- [8] S. Ross and D. Bagnell, “Efficient reductions for imitation learning,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 661–668, JMLR Workshop and Conference Proceedings, 2010.
- [9] S. Ross, G. Gordon, and D. Bagnell, “A reduction of imitation learning and structured prediction to no-regret online learning,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 627–635, JMLR Workshop and Conference Proceedings, 2011.
- [10] K. Brantley, W. Sun, and M. Henaff, “Disagreement-regularized imitation learning,” in *International Conference on Learning Representations*, 2019.
- [11] J. Ho and S. Ermon, “Generative adversarial imitation learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [12] T. Xu, Z. Li, and Y. Yu, “Error bounds of imitating policies and environments,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 15737–15749, 2020.

- [13] M. A. Bashiri, B. Ziebart, and X. Zhang, "Distributionally robust imitation learning," *Advances in neural information processing systems*, vol. 34, pp. 24404–24417, 2021.
- [14] A. Y. Ng and S. J. Russell, "Algorithms for inverse reinforcement learning," in *Proceedings of the Seventeenth International Conference on Machine Learning, ICML '00*, (San Francisco, CA, USA), p. 663–670, Morgan Kaufmann Publishers Inc., 2000.
- [15] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *International Conference on Machine Learning*, p. 1, 2004.
- [16] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3, AAAI'08*, p. 1433–1438, AAAI Press, 2008.
- [17] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization," in *Proceedings of the 33rd International Conference on Machine Learning - Volume 48, ICML'16*, p. 49–58, JMLR.org, 2016.
- [18] S. Arora, S. Du, S. Kakade, Y. Luo, and N. Saunshi, "Provable representation learning for imitation learning via bi-level optimization," in *International Conference on Machine Learning*, pp. 367–376, PMLR, 2020.
- [19] N. Rajaraman, Y. Han, L. F. Yang, K. Ramchandran, and J. Jiao, "Provably breaking the quadratic error compounding barrier in imitation learning, optimally," *arXiv preprint arXiv:2102.12948*, 2021.
- [20] J. Chang, M. Uehara, D. Sreenivas, R. Kidambi, and W. Sun, "Mitigating covariate shift in imitation learning via offline data with partial coverage," *Advances in Neural Information Processing Systems*, vol. 34, pp. 965–979, 2021.
- [21] Y. Xu, W. Gao, and D. Hsu, "Receding horizon inverse reinforcement learning," in *Advances in Neural Information Processing Systems* (A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, eds.), 2022.
- [22] G. N. Iyengar, "Robust dynamic programming," *Mathematics of Operations Research*, vol. 30, no. 2, pp. 257–280, 2005.
- [23] A. Nilim and L. El Ghaoui, "Robust control of Markov decision processes with uncertain transition matrices," *Operations Research*, vol. 53, no. 5, pp. 780–798, 2005.
- [24] N. Corporation, "Closing the sim2real gap with nvidia isaac sim and nvidia isaac replicator," 2021.
- [25] H. Xu and S. Mannor, "Distributionally robust Markov decision processes," in *Advances in Neural Information Processing Systems*, pp. 2505–2513, 2010.
- [26] W. Wiesemann, D. Kuhn, and B. Rustem, "Robust Markov decision processes," *Mathematics of Operations Research*, vol. 38, no. 1, pp. 153–183, 2013.
- [27] P. Yu and H. Xu, "Distributionally robust counterpart in Markov decision processes," *IEEE Transactions on Automatic Control*, vol. 61, no. 9, pp. 2538–2543, 2015.
- [28] S. Mannor, O. Mebel, and H. Xu, "Robust mdps with k-rectangular uncertainty," *Mathematics of Operations Research*, vol. 41, no. 4, pp. 1484–1509, 2016.
- [29] R. H. Russel and M. Petrik, "Beyond confidence regions: Tight bayesian ambiguity sets for robust mdps," *Advances in Neural Information Processing Systems*, 2019.
- [30] W. Yang, L. Zhang, and Z. Zhang, "Towards theoretical understandings of robust markov decision processes: Sample complexity and asymptotics," 2021.
- [31] K. Panaganti and D. Kalathil, "Sample complexity of robust reinforcement learning with a generative model," in *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, vol. 151 of *Proceedings of Machine Learning Research*, pp. 9582–9602, PMLR, 28–30 Mar 2022.
- [32] Z. Xu*, K. Panaganti*, and D. Kalathil, "Improved sample complexity bounds for distributionally robust reinforcement learning," Conference on Artificial Intelligence and Statistics, 2023.
- [33] K. Panaganti and D. Kalathil, "Robust reinforcement learning using least squares policy iteration with provable performance guarantees," in *International Conference on Machine Learning (ICML)*, pp. 511–520, 2021.
- [34] K. Panaganti, Z. Xu, D. Kalathil, and M. Ghavamzadeh, "Robust reinforcement learning using offline data," *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [35] Y. Wang and S. Zou, "Online robust reinforcement learning with model uncertainty," *Advances in Neural Information Processing Systems*, vol. 34, pp. 7193–7206, 2021.
- [36] J. Duchi and H. Namkoong, "Learning models with uniform performance via distributionally robust optimization," *arXiv preprint arXiv:1810.08750*, 2018.
- [37] A. Shapiro, "Distributionally robust stochastic programming," *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2258–2275, 2017.
- [38] R. Gao and A. Kleywegt, "Distributionally robust stochastic optimization with wasserstein distance," *Mathematics of Operations Research*, 2022.
- [39] D. Bertsimas, V. Gupta, and N. Kallus, "Data-driven robust optimization," *Math. Program.*, vol. 167, p. 235–292, feb 2018.
- [40] H. Namkoong and J. C. Duchi, "Stochastic gradient methods for distributionally robust optimization with f-divergences," *Advances in neural information processing systems*, vol. 29, 2016.
- [41] J. Blanchet, Y. Kang, and K. Murthy, "Robust wasserstein profile inference and applications to machine learning," *Journal of Applied Probability*, vol. 56, no. 3, p. 830–857, 2019.
- [42] B. Eysenbach and S. Levine, "Maximum entropy rl (provably) solves some robust rl problems," in *International Conference on Learning Representations*, 2022.
- [43] S. Meinecke, L. Thurner, and M. Braun, "Review of steady-state electric power distribution system datasets," *Energies*, vol. 13, no. 18, 2020.
- [44] M. Travers, "10 best healthcare data sets & examples," 2021.
- [45] S. Ettinger, S. Cheng, B. Caine, C. Liu, H. Zhao, S. Pradhan, Y. Chai, B. Sapp, C. Qi, Y. Zhou, Z. Yang, A. Chouard, P. Sun, J. Ngiam, V. Vasudevan, A. McCauley, J. Shlens, and D. Anguelov, "Large scale interactive motion forecasting for autonomous driving : The waymo open motion dataset," *arXiv*, 2021.
- [46] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.
- [47] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033, 2012.
- [48] N. Rajaraman, L. Yang, J. Jiao, and K. Ramchandran, "Toward the fundamental limits of imitation learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 2914–2924, 2020.
- [49] W. Yang, L. Zhang, and Z. Zhang, "Towards theoretical understandings of robust markov decision processes: Sample complexity and asymptotics," *arXiv preprint arXiv:2105.03863*, 2021.
- [50] K. Panaganti and D. Kalathil, "Sample complexity of robust reinforcement learning with a generative model," in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 9582–9602, 2022.
- [51] R. Chen, I. C. Paschalidis, et al., "Distributionally robust learning," *Foundations and Trends® in Optimization*, vol. 4, no. 1-2, pp. 1–243, 2020.
- [52] T. M. Cover and J. A. Thomas, "Information theory and the stock market," *Elements of Information Theory*. Wiley Inc., New York, pp. 543–556, 1991.
- [53] A. L. Gibbs and F. E. Su, "On choosing and bounding probability metrics," *International statistical review*, vol. 70, no. 3, pp. 419–435, 2002.
- [54] A. Basu, H. Shioya, and C. Park, *Statistical Inference: The Minimum Distance Approach*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability, CRC Press, 2011.
- [55] D. Bertsekas, *Reinforcement learning and optimal control*. Athena Scientific, 2019.
- [56] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun, "Reinforcement learning: Theory and algorithms," *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 2019.
- [57] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International Conference on Machine Learning*, pp. 1582–1591, 2018.
- [58] A. Raffin, "Rl baselines3 zoo." <https://github.com/DLR-RM/rl-baselines3-zoo>, 2020.
- [59] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020.
- [60] D. Kraft, "A software package for sequential quadratic programming,"

Forschungsbericht- Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt, 1988.

- [61] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.