

Smooth Imitation Learning via Smooth Costs and Smooth Policies

Sapana Chaudhary
Texas A&M University
College Station, Texas, USA
sapanac@tamu.edu

Balaraman Ravindran
Robert Bosch Centre for Data Science and AI
Indian Institute of Technology Madras
Chennai, India
ravi@cse.iitm.ac.in

ABSTRACT

Imitation learning (IL) is a popular approach in the continuous control setting as among other reasons it circumvents the problems of reward mis-specification and exploration in reinforcement learning (RL). In IL from demonstrations, an important challenge is to obtain agent policies that are smooth with respect to the inputs. Learning through imitation a policy that is smooth as a function of a large state-action (s - a) space (typical of high dimensional continuous control environments) can be challenging. We take a first step towards tackling this issue by using smoothness inducing regularizers on *both* the policy and the cost models of adversarial imitation learning. Our regularizers work by ensuring that the cost function changes in a controlled manner as a function of s - a space; and the agent policy is well behaved with respect to the state space. We call our new smooth IL algorithm *Smooth Policy and Cost Imitation Learning* (SPaCIL, pronounced “Special”). We introduce a novel metric to quantify the smoothness of the learned policies. We demonstrate SPaCIL’s superior performance on continuous control tasks from MuJoCo. The algorithm not just outperforms the state-of-the-art IL algorithm on our proposed smoothness metric, but, enjoys added benefits of faster learning and substantially higher average return.

KEYWORDS

imitation learning, continuous control, smooth policy, regularization, deep reinforcement learning

1 INTRODUCTION

A vast majority of problems of interest, including but not limited to autonomous control and robotics, are characterized by high dimensional continuous (real-valued) state and action spaces [34]. Recently, a multitude of reinforcement learning (RL) approaches (like DDPG [28], TRPO [41], PPO [43], SAC [14], *etc*) have been proposed to solve these challenging continuous control tasks. However, these algorithms require specification of a proper cost (or reward) function for any learning to be possible [1]. To circumvent this requirement, and to guide exploration in high dimensional continuous state-action spaces, Imitation Learning (IL) using demonstrations has been studied extensively [2, 7, 12, 17, 40].

In high dimensional continuous control, an important challenge is to obtain an input-robust agent policy, *i.e.*, a policy that varies in a controlled manner with respect to changes in the input state-action pair (see Fig. 1). This view of policy smoothness is equivalent to Lipschitz-continuity of the learned agent policy. Such a smooth policy can ensure agent safety, reduce agent energy consumption, and make agent behavior predictable in desired situations [25]. From

a control-theoretic perspective, policy smoothness certifies input-output stability (*i.e.*, finite \mathcal{L}_2 gain) during both exploration and deployment [20]. In certain environments, smoothness is desired from a visual standpoint, *e.g.*, smooth autonomous camera control while recording a live basketball match [8].

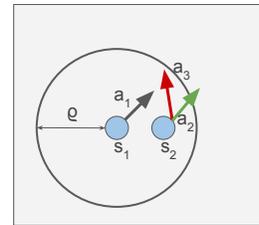


Figure 1: Understanding smoothness and non-smoothness of a policy. Consider the grey region to be a portion of the environment. We use blue circles to denote the environment states, and arrows denote the actions. Let state s_1 be perturbed to another state s_2 within a radius ϱ of s_1 . Action a_2 comes from a smooth policy as its direction and magnitude are similar to a_1 . On the other hand, a_3 comes from a non-smooth policy because of significant deviation from a_1 in both magnitude and direction.

While policy smoothness has recently gained attention in the RL community [44], it has been largely untouched in IL settings. To the best of our knowledge, policy smoothness with respect to the inputs has not been previously characterized and studied in IL. It is important to note here that by virtue of learning from demonstrations (and by employing a cost (or reward) recovery scheme), IL provides us with an additional degree of control over the final agent policy. This additional control is not available in RL. To address the challenge of obtaining smooth agent policies in high dimensional continuous environments via IL using demonstrations, we propose a smooth IL algorithm: *Smooth Policy and Cost Imitation Learning* (SPaCIL). SPaCIL learns smooth agent policies by regularizing both the cost and the policy optimization steps of adversarial imitation learning setup. Adversarial IL is a state-of-the-art IL framework that uses Inverse Reinforcement Learning (IRL) approach [30, 38], and recovers both a cost function and an agent policy. We demonstrate that our dual regularization of parameterized function approximators (for the cost and the policy) can assure the desired smoothness. At a higher level, the regularizers work by penalizing drastic changes in the cost and the policy at each step of the agent learning process. The policy regularization step is essential to control the smoothness of the learned policy model directly. The need for cost regularization stems from the observation that a cost function is a succinct definition of

a task, and by imposing proper structure on costs, we can not only recover better costs but guide our policy optimization step towards desirable policies. It is important to note that our regularizers are not tied to a specific IL objective but are general entities that can be used with wide array of algorithms. Additionally, as a byproduct of our goal, we show that we achieve considerable gains in the training and performance of the agent policy.

Related Work. A few recent works [5, 16, 19, 22, 24, 26, 29, 44, 48, 49] address some or the other notion of smoothness in adjacent/orthogonal problem settings. However, none of these works deal with smoothness in IL. The work by Kim et al. [22] on using reinforcement learning (RL) for the autonomous helicopter flight uses quadratic penalties for actions to modify the overall cost to encourage small actions and smooth control of the helicopter. That work deals with a low dimensional finite action setting and explicitly weights the cost associated with each action by a fixed factor. Such a weighing is not possible in high dimensional s - a spaces. For learning controllers for autonomous blimps, Ko et al. [24] parameterize their controller using the slope of a policy-smoothing function that determines the controller’s smoothness near a certain switching curve. However, a definite expression for policy-smoothing function exists, and the smoothness is desired over a limited region of the state space. This setting is in stark contrast to our setting, where we desire smoothness over the entire state space. Le et al. [26] study the problem of smoothly imitating an expert behaviour (using structured prediction) by ensuring that actions for adjacent states along a trajectory are similar, irrespective of proximity of adjacent states in the state space. In contrast, our notion of smoothness is not conditioned on the trajectory information and we work with the idea of smoothness of the policy space with respect to the state space. The works by Blonde et al. [5] and Shen et al. [44] are closest to our goal. Blonde et al. [5] show that enforcing Lipschitz-continuity of the learned reward function is essential for off-policy imitation learning to work well. Shen et al. [44] discuss a policy regularizer to obtain smooth policies in Reinforcement Learning (RL). We, on the other hand, focus on obtaining smooth policies in IL. Additionally, Shen et al.’s [44] work does not provide a proper evaluation of smoothness. Smoothness-inducing regularizers, similar to Shen et al.’s [44], have been previously discussed in the context of semi-supervised learning, unsupervised domain adaptation, and harnessing adversarial examples [16, 19, 29, 48, 49]. State-of-the-art IL algorithm, GAIL [17] introduces various cost function regularizers to obtain various instances of IL algorithms; however, none of those regularizes enforce any special structure on the costs.

In what follows, we describe our smooth IL approach in greater detail. The main contributions of the paper are as follows:

- (1) Formalization of the notion of smooth policies using Lipschitz continuity.
- (2) Theoretical study of how Lipschitz continuous rewards facilitate in obtaining Lipschitz continuous agent policy in on-policy continuous control.
- (3) Introduction of smoothness inducing cost function and policy function regularizers to realize a smooth IL algorithm.
- (4) Introduction of a novel metric (that captures Lipschitz continuity of a learned policy) to assess the smoothness of a learned policy.

- (5) Empirical testing of smoothness of the learned policies, and validation of other claims on high dimensional continuous control tasks.

2 BACKGROUND

Markov Decision Process. We consider gamma discounted infinite horizon continuous Markov Decision Processes (MDPs) [45] as the core framework. An MDP is specified by the tuple $\langle \mathcal{S}, \mathcal{A}, P, c, \gamma, \rho_0 \rangle$, where $\mathcal{S} \subset \mathbb{R}^{D_s}$ and $\mathcal{A} \subset \mathbb{R}^{D_a}$ are compact sets with non-zero Lebesgue measure. \mathcal{S} is the set of possible states, \mathcal{A} is the set of possible actions and $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the cost function. D_s and D_a are the dimensions of the state and the action spaces, respectively. P for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ triplet gives the probability of moving from state s to state s' on taking action a at s . γ is the discount factor and ρ_0 is the initial state distribution.

Policy, occupancy measure, and expected cost. A stationary stochastic policy, $\pi(a|s)$ gives agent’s behaviour in the environment. Return, G , is a measure of goodness of a policy and is defined as $G(\pi) = -\mathbb{E}_\pi[c(s, a)] = -\mathbb{E}[\sum_{t \geq 0} \gamma^t c(s_t, a_t)]$. The return is estimated from trajectories as: $\hat{G}(\pi) = -\mathbb{E}_\tau[\sum_{t=0}^T \gamma^t c(s_t, a_t)]$, where $\tau \sim \pi$ is a trajectory of the form $\{(s_i, a_i)\}_{i=1}^T$. Here $s_0 \sim \rho_0$ is the starting state, $a_t \sim \pi(\cdot|s_t)$ and $s_{t+1} \sim P(\cdot|s_t, a_t)$. T denotes the time step until which we sample a trajectory. There is a one-to-one mapping between a policy, π and its occupancy measure, $\rho_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defined as $\rho_\pi(s, a) = \pi(a|s) \sum_{t=0}^{\infty} \gamma^t P(s_t = s|\pi)$. Expected cost in terms of ρ_π is given by $\mathbb{E}_\pi[c(s, a)] = \int_{\mathcal{S} \times \mathcal{A}} \rho_\pi(s, a) c(s, a) ds da$. The stationary distribution of the policy π is denoted by $\rho_\pi(s)$. The expert policy is denoted by π_E . For a general function of states, $f(s)$, $\mathbb{E}_{s \sim \rho}[f(s)]$ means a γ -discounted expectation (as is for $G(\pi)$), unless stated otherwise.

Value functions and advantage. The value function V^π and action value function Q^π can be written as $V^\pi(s) = -\mathbb{E}_\pi[c(\cdot, \cdot) | s_0 = s]$ and $Q^\pi(s, a) = -\mathbb{E}_\pi[c(\cdot, \cdot) | s_0 = s, a_0 = a]$. Advantage function $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ reflects the expected additional cost that the agent bears after taking action a in state s .

Parameterized representations. We use parameterized function approximators (neural networks) for realization of the cost and the policy. When needed, we denote the cost function as $c_\omega(s, a)$ and the policy as $\pi_\theta(a|s)$, where $\omega \in \Omega \subset \mathbb{R}^{D_\omega}$ and $\theta \in \Theta \subset \mathbb{R}^{D_\theta}$. D_ω and D_θ are the dimensions of the parameters. We assume that Ω and Θ are compact sets. \mathcal{C} is the space of all cost functions. Π is the space of all policies. $\|\cdot\|$ denotes the L^2 norm unless specified otherwise. $|\cdot|$ represents the usual norm.

Imitation Learning. IL solution approaches can be broadly divided into: Behaviour Cloning (BC) [4, 39] and imitation learning using inverse reinforcement learning (IRL, [1, 11, 27, 50]). For high dimensional continuous control tasks, IL using IRL is the method of choice as BC suffers from covariate shift errors [35, 37]. IL using IRL can be cast as the following bi-level optimization problem (GAIL, [17]):

$$\text{IL}(\pi_E) = \arg \max_{c \in \mathcal{C}} -\psi(c) + \left(\min_{\pi \in \Pi} \mathbb{E}_\pi[c(s, a)] \right) - \mathbb{E}_{\pi_E}[c(s, a)], \quad (1)$$

where $\psi : C \rightarrow \overline{\mathbb{R}}$ is a closed proper convex cost function loss specifier, *i.e.*, it helps specify a trainable loss for the cost function model. In the first level of optimization, we fix the cost function and update the policy using $\min_{\pi \in \Pi} \mathbb{E}_{\pi} [c(s, a)]$. Next, we fix the policy from the previous update, and update the cost function using $\arg \max_{c \in C} -\psi(c) + \mathbb{E}_{\pi} [c(s, a)] - \mathbb{E}_{\pi_E} [c(s, a)]$.

Policy optimization. The policy update step (performed using the trust region based policy optimization algorithm (TRPO, [41])) is given by:

$$\pi_{k+1} = \arg \min_{\pi} - \mathbb{E}_{s \sim \rho^{\pi_k}, a \sim \pi_k} \left[\frac{\pi(a | s)}{\pi_k(a | s)} A^{\pi_k}(s, a) \right]$$

$$\text{subject to } \mathbb{E}_{s \sim \rho^{\pi_k}} [\mathcal{D}_{\text{KL}}(\pi_k(\cdot | s) \| \pi(\cdot | s))] \leq \delta. \quad (2)$$

Here, $\mathcal{D}_{\text{KL}}(\pi_k(\cdot | s) \| \pi(\cdot | s)) = \mathbb{E}_{a \sim \pi_k} \left[\frac{\pi_k(a | s)}{\pi(a | s)} \right]$ is the Kullback–Leibler (KL) divergence and the constraint (in Eqn. 2) bounds the KL-divergence between two consecutive policies by δ . Here, the expectation in constraint is not γ -discounted.

Cost optimization. Using the cost function loss specifier $\psi_{\text{GA}}(c)$ of GAIL [17] the cost parameters are updated as:

$$c_{k+1} = \arg \max_c \mathbb{E}_{s \sim \rho^{\pi_k}} [\log(D(s, \pi_k(\cdot | s)))] +$$

$$\mathbb{E}_{s \sim \rho^{\pi_E}} [\log(1 - D(s, \pi_E(\cdot | s)))] \quad (3)$$

where $\log(D(s, \pi_k(\cdot | s))) = c(s, \pi_k(\cdot | s))$ is the cost function and $D : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1)$ is the classifier.

Gaussian policy representation. Additionally, to include stochasticity [31] in the on-policy [45] policy gradient approach (TRPO) we consider our stationary stochastic policies to be Gaussian distributed with $\mu_{\theta}(s)$ as the mean and standard deviation (std) given by a fixed quantity σ . Here, $\mu_{\theta}(s)$ is a deterministic function of the states parameterized by θ . Thus, an action a sampled from our policy can be written as,

$$a \sim \mathcal{N}(\mu_{\theta}(s), \sigma) \implies a = \mu_{\theta}(s) + \sigma z, \quad z \sim \mathcal{N}(0, 1), \quad (4)$$

where $\mathcal{N}(\cdot)$ is the standard Normal distribution.

3 PROBLEM DEFINITION

3.1 Defining smooth policy, smooth cost, and smooth transition model

This section formally defines a smooth policy, a smooth cost, and a smooth transition model.

DEFINITION 1 (SMOOTH POLICY). *Let ξ be a metric on the space of policies, Π . A stationary stochastic policy $\pi(a|s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is smooth with respect to the inputs, \mathcal{S} , if for all $s_1, s_2 \in \mathcal{S}$ it is Lipschitz continuous and hence, there exist an $M_{\pi} \geq 0$ such that*

$$\forall s_1, s_2 \in \mathcal{S}, \quad \xi(\pi(\cdot | s_1), \pi(\cdot | s_2)) \leq M_{\pi} \|s_1 - s_2\|. \quad (5)$$

DEFINITION 2 (SMOOTH COST). *A cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is smooth with respect to the inputs, $\mathcal{S} \times \mathcal{A}$, if for all $(s_1, a_1), (s_2, a_2) \in \mathcal{S} \times \mathcal{A}$ it is Lipschitz continuous and hence, there exist an $M_c \geq 0$ such that*

$$\forall (s_1, a_1), (s_2, a_2) \in \mathcal{S},$$

$$\|c(s_1, a_1) - c(s_2, a_2)\| \leq M_c (\|s_2 - s_1\| + \|a_2 - a_1\|). \quad (6)$$

DEFINITION 3 (SMOOTH TRANSITION MODEL). *If the transition model, P is L_P -Lipschitz continuous it satisfies the following constraint for every two state action pairs $(s_1, a_1), (s_2, a_2) \in \mathcal{S} \times \mathcal{A}$ and all 1-Lipschitz value functions V :*

$$\left| \int_{s' \in \mathcal{S}} (P(s' | s_1, a_1) - P(s' | s_2, a_2)) V(s') ds' \right|$$

$$\leq L_P (\|s_1 - s_2\| + \|a_1 - a_2\|). \quad (7)$$

If V is L_v Lipschitz, the right hand side of the inequality in Eqn. 8 will be scaled by L_v .

3.2 IL with smooth policy

Smooth policies are crucial in high dimensional continuous control for diverse reasons ranging from critical (robot safety) to aesthetic (visual appeal). When a cost function can be appropriately constructed and specified by a problem designer, reinforcement learning is the go-to solution approach. However, cost function design is a tedious task, and cost functions tend to be grossly mis-specified [1, 9]. Imitation learning helps overcome the challenge of cost function design by defining strategies that learn agent policies from (expert) demonstrations. However, existing IL approaches do not guarantee that the learned policies are smooth, *i.e.*, there is no existing approach that solves the following (general) problem:

$$\pi_{\text{smooth}}^{\star} = \arg \min_{\pi \in \Pi} \mathbb{E}_{\pi} [c(s, a) | \{(s_E, a_E)\} \sim \pi_E]$$

$$\text{subject to } \xi(\pi, \pi_{\text{perturbed}}) \leq \epsilon', \quad (9)$$

where $\mathbb{E}_{\pi} [c(s, a) | \{(s_E, a_E)\} \sim \pi_E]$ represents a general IL objective, $\pi_{\text{perturbed}}$ represents a policy obtained by perturbing the original policy (π) by a small amount, and ϵ' captures desired policy smoothness. The policy perturbation can be achieved in numerous ways. For our purpose of policy smoothness, $\pi_{\text{perturbed}}$ is obtained by inducing controlled amount of noise in the states along trajectories sampled by a policy π . Hence, the problem is to obtain optimal smooth agent policy, $\pi_{\text{smooth}}^{\star}$ from the class of policies, Π using expert demonstration data (N trajectories of the form: $\tau_i^E = \{(s_j, a_j)\}_{j=1}^T, i = 1, 2, \dots, N$). The policy optimization is subject to the constraint: $\xi(\pi, \pi_{\text{perturbed}}) \leq \epsilon'$, *i.e.*, choose that policy which when perturbed in a controlled manner behaves similar to the original unperturbed policy. The similarity in behaviour is characterized using the distance metric ξ defined over Π .

4 METHOD FOR OBTAINING SMOOTH POLICIES IN IL

We take the approach of adversarial imitation learning [12, 17] to solve the problem discussed above (Eqn. 9). Adversarial IL is a high dimensional counterpart of IL using IRL, and is mathematically formulated as:

$$\text{IL}(\pi_E) = \arg \max_{c \in C'} \left(\min_{\pi \in \Pi} \mathbb{E}_{\pi} [c(s, a)] \right) - \mathbb{E}_{\pi_E} [c(s, a)], \quad (10)$$

where C' is the cost function class modified by the loss function chosen to train the cost model. In this approach, the agent policy (π) and the cost function (c) are simultaneously learned using bi-level optimization of parameterized models. The cost function is learned using the expert demonstration data and the samples from the current

agent policy model. The cost model update rule is defined in Eqn. 3. The agent policy is learned using trust region-based policy optimization algorithm (Eqn. 2). To achieve the goal of policy smoothness, we propose to include smoothness inducing regularizers in both the policy and the cost optimization steps of adversarial imitation learning:

$$\text{IL}(\pi_E) = \arg \max_{c \in \mathcal{C}'} \left(\min_{\pi \in \Pi} \mathbb{E}_{\pi} [c(s, a)] + \mathcal{R}_{(s)}^{\pi} \right) - \mathbb{E}_{\pi_E} [c(s, a)] - \mathcal{R}_{(s, a)}^c, \quad (11)$$

where $\mathcal{R}_{(s)}^{\pi}$ is the smoothness inducing regularizer on policy, and $\mathcal{R}_{(s, a)}^c$ is the smoothness inducing regularizer on cost function. The exact forms of these regularizers is discussed shortly. While policy regularization using a regularizer that captures the constraint in Eqn. 9 seems like an obvious first approach, it is not immediately clear as to why we need to regularize the cost function. In the following section we describe the need for cost smoothness.

4.1 Smooth costs assist in learning a smooth policy

A cost function is a central tool to perform learning through interaction. Using the following theorems, we show that a smooth cost function helps in obtaining smooth RL policies (, and hence in IL using IRL). We show that a Lipschitz continuous (smooth) cost function ensures resulting optimal value functions $V^*(s)$ and $Q^*(s, a)$ are Lipschitz continuous (smooth). We then show that if a Lipschitz continuous mapping is used to obtain π^* from Q^* , then π^* is Lipschitz continuous. These results are of independent interest apart from providing a clear motivation to regularize cost function learning in our method.

THEOREM 1 (GENERALIZATION OF THEOREM 5.9 IN [32] FOR THE CASE OF CONTINUOUS STATE AND ACTION SPACE). *For a given MDP, if the cost function, $c(s, a)$ and the transition model, $P(s'|s, a)$ are L_c and L_p -Lipschitz continuous, respectively, with respect to the state and the action pairs and $\gamma L_p < 1$ then*

- the optimal state value function $V^*(s)$ is $\frac{L_c}{1-\gamma L_p}$ Lipschitz with respect to states, and*
- the optimal state action value function $Q^*(s, a)$ is $\frac{L_c}{1-\gamma L_p}$ Lipschitz continuous with respect to states and actions.*

This theorem is a generalization of Theorem 5.9 in [32] for the case of continuous state and action space. The proof follows along the lines of one in [32] by replacing discrete Bellman optimality operator for \mathcal{V} and \mathcal{Q} with continuous counterparts. The proof is included in Appendix A.1.

THEOREM 2. *Let κ be a pseudo-metric on the space of conditional state-action value functions $Q(s, \cdot)$. Let $H : Q(s, \cdot) \rightarrow \mu(s)$ be a smooth mapping that outputs (near) greedy stationary mean policies with respect to the input conditional state action value function $Q(s, \cdot)$ such that*

$$\|H(Q(s_1)) - H(Q(s_2))\| \leq L_{\mu} \kappa(Q(s_1, \cdot), Q(s_2, \cdot)), \forall s_1, s_2 \in \mathcal{S}.$$

For a given $\frac{L_c}{1-\gamma L_p}$ -Lipschitz continuous state action value function $Q^(s, a)$, the policy obtained as $\mu^*(s) = H(Q^*(s, \cdot))$ is Lipschitz continuous with respect to states. Then, the stationary stochastic policy π^* obtained from μ^* as $\mathcal{N}(\mu^*, \sigma)$ (for a fixed σ) is Lipschitz continuous with respect to states.*

The proof is included in Appendix A.2.

4.2 Exact forms of regularizers

With the said motivations and goal in mind, we now discuss our regularizers' exact forms to achieve *smooth* policies in IL.

Policy Smoothing. The aim of the policy regularizer is to encourage the agent policy to be smooth within an ϵ neighbourhood of all the states sampled according to the current policy model. At each iteration of the policy update step, we sample N trajectories from the current policy where each trajectory is of the form $\tau_i = \{(s_i^t, a_i^t)\}_{t=1}^T$. We then nudge every state, s , in the sampled batch to obtain, \tilde{s} , such that \tilde{s} lies in an ϵ -radius ball around s (i.e., $\tilde{s} \in \mathbb{B}_d(s, \epsilon)$) and the policy, π_{θ} changes the most at this s' . The maximum policy change is defined using a suitable divergence $\mathcal{D}(\pi_{\theta}(s), \pi_{\theta}(\tilde{s}))$. In this work, we consider \mathcal{D} to be Jeffrey's Divergence: $\mathcal{D}_J(P||Q) = \frac{1}{2} \mathcal{D}_{\text{KL}}(P||Q) + \frac{1}{2} \mathcal{D}_{\text{KL}}(Q||P)$. This policy regularizer is similar in spirit to the one considered in [44]. At a fundamental level, the regularizer measures the local Lipschitz continuity of policy π_{θ} under the divergence \mathcal{D} and thus limits the policy output decision change if a small perturbation is added to a certain state s .

$$\mathcal{R}_S^{\pi}(\theta_k) = \mathbb{E}_{s \sim \rho^{\pi_{\theta_k}}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{D}_J(\pi_{\theta_k}(\cdot | s), \pi_{\theta_k}(\cdot | \tilde{s})). \quad (12)$$

Note that the term inside expectation is γ discounted. This regularizer is now added to the policy optimization (TRPO) objective function to obtain smooth policy update rule:

$$\begin{aligned} \theta_{k+1} = \arg \min_{\theta} & - \mathbb{E}_{s \sim \rho_{\theta_k}^{\pi}, a \sim \pi_{\theta_k}} \left[\frac{\pi_{\theta}(a | s)}{\pi_{\theta_k}(a | s)} A^{\pi_{\theta_k}}(s, a) \right] + \\ & \lambda_1 \mathbb{E}_{s \sim \rho_{\theta_k}^{\pi}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{D}_J(\pi_{\theta}(\cdot | s) || \pi_{\theta}(\cdot | \tilde{s})), \quad (13) \\ & \text{subject to } \mathbb{E}_{s \sim \rho^{\pi_{\theta_k}}} [\mathcal{D}_{\text{KL}}(\pi_{\theta_k}(\cdot | s) || \pi_{\theta}(\cdot | s))] \leq \delta. \end{aligned}$$

Cost Smoothing. To obtain smooth costs, we propose to regularize the cost optimization step. The goal of this regularizer is to regulate the worst change in cost function within an ϵ -ball around states obtained from the mixture trajectories, $\hat{\tau} \leftarrow \zeta \tau_E + (1 - \zeta) \tau_I$, where ζ is a mixing parameter. The smoothness is induced for (s, a) pairs sampled from the mixture trajectory to generalize smoothness across the entire state-action space.

A few more words about trajectory mixing are warranted here. The mixing of trajectories considered in our work is reminiscent of policy mixing introduced in RL in Kakade et al. [21] and in IL introduced in Ross et al. [36]. Policy mixing in these works aims to ensure conservative policy update at each policy iteration step. A greedy policy update based purely on *approximate* state-action values has shown to affect the policy learning process adversely. A similar mixing of expert and non-expert data is used in the regularizer sampling distribution of WGAN-GP [13]. The goal of data mixing in WGAN-GP is to ensure Lipschitzness property is satisfied by the gradients (of a model) for both the expert and the non-expert data. Our regularizer draws inspiration from the observations mentioned above. We mix agent and expert trajectories to ensure costs are smoothed both as a function of agent and expert s - a pairs and not just imperfect agent data. Such mixing guarantees conservative enforcement of smoothness regularizer over cost function, rather than changing cost function drastically over iterations, making the learning algorithm unstable. Additionally, our regularizer provides a

consistent learning signal to the cost model even when the agent and the expert policy supports do not overlap [3, 13]. The cost regularizer takes the following exact form:

$$\mathcal{R}_{(s,a)}^c(\theta) = \mathbb{E}_{s \sim \rho^{\tilde{\pi}_\theta}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \|c(s, \pi_\theta(\cdot | s)) - c(s, \pi_\theta(\cdot | \tilde{s}))\|_2. \quad (14)$$

Here, $\tilde{\pi}_\theta$ is a mixture policy whose exact form is not required for our purposes. We only need samples, τ from this policy (which are obtained as $\hat{\tau} \leftarrow \zeta \tau_E + (1 - \zeta) \tau_i$, where ζ is the mixing parameter). Using the regularizer of Eqn. 14, the cost optimization problem becomes:

$$\begin{aligned} w_{k+1} = & \arg \max_w \mathbb{E}_{s \sim \rho^{\pi_{\theta_k}}} [\log(D_w(s, \pi_{\theta_k}(\cdot | s)))] + \\ & \mathbb{E}_{s \sim \rho^{\pi_E}} [\log(1 - D_w(s, \pi_{\theta_E}(\cdot | s)))] - \\ & \lambda_2 \mathbb{E}_{s \sim \rho^{\tilde{\pi}_\theta}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \|c(s, \pi_\theta(\cdot | s)) - c(s, \pi_\theta(\cdot | \tilde{s}))\|_2. \end{aligned} \quad (15)$$

5 PRACTICAL ALGORITHM

Maximization over ϵ -Ball. Both the regularized optimization problems in Eqn. 13 and Eqn. 15 require us to solve a maximization problem over the ϵ -ball around a certain state s . The goal of this maximization problem is to find a state s' within an ϵ -ball of a state s at which a certain function, $f(s, s')$ takes the maximum value. In Algorithm 1, we discuss a general projected gradient based approach to solve this maximization problem that is applicable to both Eqn. 13 and Eqn. 15. For the policy smoothing regularizer of Eqn. 12 the f in Algorithm 1 is given by $f(s, s') = D_J(\pi_\theta(s), \pi_\theta(s'))$ (the Jeffrey’s divergence between policies $\pi_\theta(s)$ and $\pi_\theta(s')$). For the cost smoothing regularizer of Eqn. 14, $f(s, s') = \|c(s, \pi_\theta(\cdot | s)) - c(s, \pi_\theta(\cdot | s'))\|_2$ (the L_2 distance between costs). Now that we have a procedure to obtain the regularizers, our proposed Smooth IL algorithm is provided in Algorithm 2.

Algorithm 1 Maximization of f over ϵ -ball of a state s

- 1: **Input:** s, ϵ, η_δ
 - 2: **Initialize:** δ_0
 - 3: (N steps of projected gradient descent)
 - 4: **for** $\ell = 1, 2, \dots, N-1$
 - 5: **(Update δ in the direction of increase in f)**
 - 6: $\delta_{\ell+1} = \delta_\ell + \eta_\delta \nabla_\delta f(s, s + \delta_\ell)$
 - 7: **(Project $\delta_{\ell+1}$ onto the ϵ -ball)**
 - 8: $\delta_{\ell+1} = \Pi_{\mathbb{B}_d(0, \epsilon)}(\delta_{\ell+1})$
 - 9: **end for**
-

5.1 Evaluating smoothness of the learned policy

To quantify the smoothness of the learned agent policy, we introduce the following (general) novel metric that captures local Lipschitz continuity of the policy:

$$J(\pi) = \mathbb{E}_{s \sim \rho^\pi} \left[\max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{D}_J(\pi(\cdot | s), \pi(\cdot | \tilde{s})) \right], \quad (16)$$

where the term inside expectation is *not* γ -discounted. We do not include discounting here because we desire policy at any state sampled by the policy to be smooth. For the particular case of Gaussian policies, to assess the Lipschitz continuity of a stochastic policy π , we can look at the Lipschitz smoothness of its deterministic mean

Algorithm 2 SPaCIL: Smooth Reward and Policy Imitation Learning

- 1: **Input:** Expert trajectories $\tau_E \sim \pi_E$, initial policy and discriminator parameters θ_0, w_0
 - 2: **for** $k = 1, 2, \dots, K$
 - 3: **Smooth policy update:**
 - 4: Sample N trajectories $\tau_i^k \sim \pi_{\theta_k}, i = 1, 2, \dots, N$
 - 5: Estimate regularizer in Eqn. 12 using Algorithm 1
 - 6: Update policy using regularized TRPO of Eqn. 13
 - 7: **Smooth cost update:**
 - 8: Estimate regularizer in Eqn. 14 using Algorithm 1
 - 9: Update cost using Eqn. 15
 - 10: **end for**
-

function μ (see Eqn. 4). For this case, practical smoothness metric takes the following form:

$$J(\{\tau_i\}_{i=1}^N) = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^T \max_{\tilde{s}_j \in \mathbb{B}_d(s_j, \epsilon)} \frac{\|\mu(s_j) - \mu(\tilde{s}_j)\|}{\|s - \tilde{s}_j\|}, s_j \sim \tau_i. \quad (17)$$

We do not (on purpose) consider a global Lipschitz constant L corresponding to the deterministic mean function, $\mu(s)$ of the following form:

$$L = \max_{s_1 \neq s_2} \frac{\|\mu(s_1) - \mu(s_2)\|}{\|s_1 - s_2\|}, \quad s_1, s_2 \in \mathcal{S}. \quad (18)$$

The choice to quantify smoothness using local Lipschitz constant stems from the fact that finding the maximum over \mathcal{S} in Eqn. 18 is impractical for high dimensional environments.

6 EXPERIMENTS

Overview. In this section, we aim at investigating the following research questions (RQs):

- (1) Using average return (G) as a metric, how well does the learned agent policy for SPaCIL and GAIL behave in the environment? (*Sec. 6.1, Table 1, Fig. 2*)
- (2) Does regularization of the policy and the cost space result in faster learning for SPaCIL? (*Sec. 6.1, Fig. 2*)
- (3) Using our smoothness metric (J), which algorithm gives the best resulting smooth policies? (*Sec. 6.1, Table 2*)
- (4) How good is our smoothness evaluation metric (J)? If we sufficiently perturb our smooth agent policy model, does J worsen? (*Sec. 6.2, Fig. 3*)
- (5) How do the two regularizations affect SPaCIL’s performance and in what ways? (*Sec. 6.3, Fig. 4*)

We answer these questions by performing multiple experiments on continuous control tasks from MuJoCo [47]. We specifically work with these environments: Reacher-v2, Hopper-v2, Walker2d-v2, HalfCheetah-v2, and Ant-v2. We would drop v2 from these environments in all future references to them.

Implementation details. For all the environments, the expert is trained using TRPO [41]. We then sample trajectories from the best expert model and form our demonstration dataset. The demonstration dataset does not necessarily come from a smooth expert policy because we do not explicitly incorporate smoothness during TRPO

Table 1: Average Return (and one standard deviation) from 100 trajectories sampled from 5 different best agent models (i.e., five different random seeds). We observe that SPaCIL outperforms GAIL by a substantial margin, and enjoys much lesser variance in average return across different runs.

	Reacher	Hopper	Walker2d	HalfCheetah	Ant
S / \mathcal{A}	$\mathbb{R}^{11} / \mathbb{R}^2$	$\mathbb{R}^{11} / \mathbb{R}^3$	$\mathbb{R}^{17} / \mathbb{R}^6$	$\mathbb{R}^{17} / \mathbb{R}^6$	$\mathbb{R}^{111} / \mathbb{R}^8$
Expert (non-smooth)	-4.18 ± 1.79	3562.04 ± 26.61	4224.34 ± 18.79	4022.09 ± 79.72	4872.84 ± 568.69
GAIL	-5.27 ± 2.72	3326.43 ± 872.42	4275.00 ± 14.69	4064.12 ± 120.30	4587.81 ± 117.84
SPaCIL	-4.39 ± 1.61	3662.71 ± 39.27	4397.58 ± 6.92	4141.10 ± 93.84	4788.73 ± 75.18

Table 2: Smoothness metric (J) (and one standard deviation) of the best trained models estimated over 150 randomly sampled trajectories from five different policy models. SPaCIL learns a substantially smooth policy irrespective of how smoothness encoded in expert demonstrations.

	Reacher	Hopper	Walker2d	HalfCheetah	Ant
Expert (non-smooth)	$1.77e-05 \pm 3.34e-06$	12.34 ± 0.089	56.61 ± 0.031	23.88 ± 11.70	$6.4 \pm 3.81e-3$
GAIL	$5.86e-5 \pm 9.79e-5$	12.91 ± 0.85	117.31 ± 6.65	24.25 ± 0.77	1.543 ± 0.62
SPaCIL	$2.69e-5 \pm 2.02e-5$	7.77 ± 0.067	35.61 ± 0.47	17.73 ± 0.12	0.93 ± 0.09

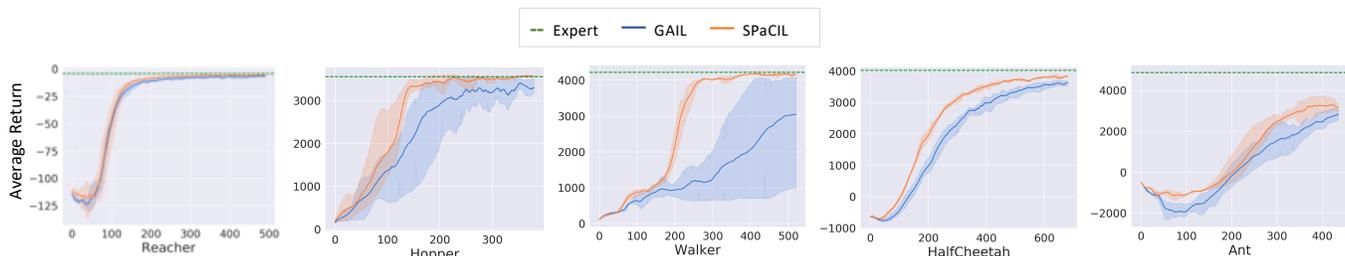


Figure 2: Learning curves for GAIL and SPaCIL on challenging continuous control tasks. The x axis is the number of algorithmic iterations. The dotted green line represents the demonstration dataset’s average return. The dual regularization of SPaCIL results in faster learning.

training. This setting is practical as in reality we might have non-smooth demonstrations but would still desire a smooth controller. The imitation learning algorithms are trained on these datasets. The algorithms do not have access to the actual environment rewards and dynamics. To make the comparison fairer, for each environment, all the hyperparameters of GAIL and SPaCIL are the same except for SPaCIL’s regularization-specific hyperparameters. More implementation details can be found in Appendix B.

6.1 Policy average return and smoothness

As the first set of experiments, we train GAIL and SPaCIL on demonstration datasets from each environment. The algorithms are trained for 400 to 500 iterations depending on when the average return stabilizes. The learning curves are included in Fig. 2. The average return from SPaCIL is better than GAIL across all the environments. Higher returns mean that our smoothness regularizations result in overall better-behaved agents (RQ 1). SPaCIL also enjoys lesser variance in average return, implying the learned policies will give the

said average returns with greater confidence. SPaCIL additionally performs better than the average return of the demonstration data for Hopper, Walker, and HalfCheetah.

The smoothness metric J is estimated from 150 trajectories (sampled using different random seeds) from five best agent models. Here as well, SPaCIL outperforms GAIL and recovers much smoother learned agent policies (RQ 3). The training curves and J for Walker are particularly interesting. GAIL’s inability to tap onto the smoothness (seen from very high J for this task) results in a much higher variance of the training curve. From the nature of the learning curve for Reacher, we see that the advantage of SPaCIL over GAIL is more pronounced in high-dimensional tasks. From Fig. 2 it is evident that SPaCIL results in faster stabilized returns (RQ 2).

6.2 Validating smoothness metric

This section aims to answer if the proposed smoothness metric (J) is truly meaningful. To answer this, we take the best SPaCIL policy model for HalfCheetah and Hopper, and perturb this model by adding

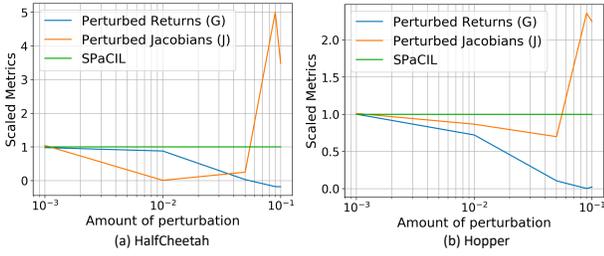


Figure 3: Variation of scaled average policy smoothness metric (J) and scaled average return for perturbed SPaCIL models for HalfCheetah and Hopper. Scaling is performed w.r.t. SPaCIL J and G , respectively. The x-axis quantifies perturbation as the standard deviation of noise added to model parameters. For sufficiently perturbed models, J is very high, validating J as a smoothness metric.

a zero mean, fixed standard deviation Gaussian noise to the model parameters. Standard deviations considered are: 0.001, 0.01, 0.005, 0.009 and 0.1. We then obtain the average return, and the average smoothness metric (J) from all these perturbed models. Fig. 3 depicts scaled average returns and average J for HalfCheetah and Hopper. The scaling is done with respect to SPaCIL model’s average return and average J (in green colour), *i.e.*, in Fig. 3 SPaCIL model’s both average return and average J are equal to one. We observe that the average return decreases by a small amount with small perturbation with severe decrease when perturbation standard deviation is 0.1 (high). This means that our learned policy model is quite robust to model parameter perturbation. The smoothness metric initially decreases showing existence of a policy that is smoother than the baseline at the cost of decrease in return. This highlights an important facet of our method - we want high performing policies to be smooth, and not desire excessive smoothness at the cost of lesser return. Also, an extremely smooth policy might mean the agent can hardly move and hence gets a low return. When the model is sufficiently perturbed (with a perturbation standard deviation of 0.1), the resulting J is very high, showing that the policy model is non-smooth (RQ 4).

6.3 Policy regularization vs. cost regularization: Which is more important?

We run Hopper with varying amounts of cost and policy regularization. Fig. 4 depicts the results from this experiment. We get the best smooth policy for $\lambda_1 = 0.001$ and $\lambda_2 = 1$ with $G = 3642.76$ and $J = 9.556$. The agent obtains the maximum return of 3715.08 at $\lambda_1 = 0$ and $\lambda_2 = 1$, *i.e.*, merely smoothing the costs results in policies that enjoy higher return. However, this policy is highly non-smooth with a $J = 68.67$. Thus, the agent has to pay some additional cost to obtain visual smoothness. We can see from Fig. 4 that a delicate balance of both policy and cost regularization is needed to obtain high performing smooth policies.

7 CONCLUSION

In this work, we develop a novel model-free on-policy IL algorithm: *Smooth Policy and Cost Imitation Learning* (SPaCIL) that learns a

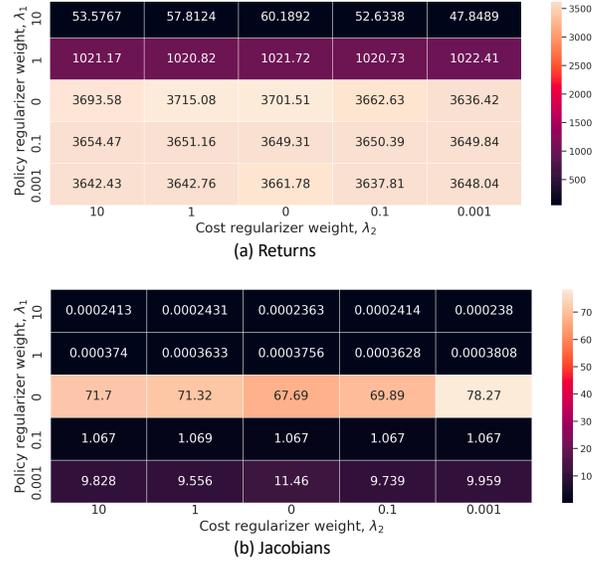


Figure 4: Average return (G) and average smoothness metric (J) variation with varying policy and cost regularization parameters (λ_1 and λ_2) for Hopper. Both policy and cost regularizations are needed to obtain high performing smooth policies.

smooth agent policy. Smoothness in policies is achieved by regularizing both the cost and the policy models of adversarial imitation learning framework. Through smoothness-inducing regularization, our algorithm can encode the domain knowledge about the smoothness of costs and policies. Our algorithm not only obtains smooth IL policies (measured by the policy smoothness metric that we introduce), it results in policies with a higher return than state-of-the-art adversarial IL algorithm GAIL. Our algorithm enjoys added benefits of faster learning and a much lower variance in average returns and smoothness metric.

REFERENCES

- [1] Pieter Abbeel and Andrew Y Ng. 2004. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*. ACM, 1.
- [2] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. 2009. A survey of robot learning from demonstration. *Robotics and autonomous systems* 57, 5 (2009), 469–483.
- [3] Martin Arjovsky and Léon Bottou. 2017. Towards principled methods for training generative adversarial networks. In *NIPS 2016 Workshop on Adversarial Training*. In review for ICLR, Vol. 2016.
- [4] Michael Bain and Claude Sammut. 1995. A Framework for Behavioural Cloning.. In *Machine Intelligence 15*. 103–129.
- [5] Lionel Blondé, Pablo Strasser, and Alexandros Kalousis. 2020. Lipschitzness Is All You Need To Tame Off-policy Generative Adversarial Imitation Learning. *arXiv preprint arXiv:2006.16785* (2020).
- [6] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. Openai gym. *arXiv preprint arXiv:1606.01540* (2016).
- [7] Sylvain Calinon. 2009. *Robot programming by demonstration*. EPFL Press.
- [8] Jianhui Chen, Hoang M Le, Peter Carr, Yisong Yue, and James J Little. 2016. Learning online smooth predictors for realtime camera planning using recurrent decision trees. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4688–4696.
- [9] Peter Dayan and Bernard W Balleine. 2002. Reward, motivation, and reinforcement learning. *Neuron* 36, 2 (2002), 285–298.
- [10] Kenneth Eriksson, Donald Estep, and Claes Johnson. 2003. *Applied mathematics: Body and soul: Calculus in several dimensions*. Vol. 3. Springer Science &

- Business Media.
- [11] Chelsea Finn, Sergey Levine, and Pieter Abbeel. 2016. Guided cost learning: Deep inverse optimal control via policy optimization. In *Proceedings of the 33rd International Conference on Machine Learning*, Vol. 48.
 - [12] Justin Fu, Katie Luo, and Sergey Levine. 2017. Learning Robust Rewards with Adversarial Inverse Reinforcement Learning. *arXiv preprint arXiv:1710.11248* (2017).
 - [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. Improved Training of Wasserstein GANs. *arXiv preprint arXiv:1704.00028* (2017).
 - [14] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, PMLR, 1861–1870.
 - [15] Robert Hecht-Nielsen. 1992. Theory of the backpropagation neural network. In *Neural networks for perception*. Elsevier, 65–93.
 - [16] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. 2019. Using self-supervised learning can improve model robustness and uncertainty. In *Advances in Neural Information Processing Systems*. 15663–15674.
 - [17] Jonathan Ho and Stefano Ermon. 2016. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*. 4565–4573.
 - [18] Martin Sleziak (<https://math.stackexchange.com/users/8297/martin-sleziak>). [n.d.]. Is the maximum function Lipschitz continuous? Mathematics Stack Exchange. <https://math.stackexchange.com/q/1742410> <https://math.stackexchange.com/q/1742410> URL:<https://math.stackexchange.com/q/1742410> (version: 2016-04-14).
 - [19] Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. 2019. Smart: Robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization. *arXiv preprint arXiv:1911.03437* (2019).
 - [20] Ming Jin and Javad Lavaei. 2018. Control-theoretic analysis of smoothness for stability-certified reinforcement learning. In *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 6840–6847.
 - [21] Sham Kakade and John Langford. 2002. Approximately optimal approximate reinforcement learning. In *In Proc. 19th International Conference on Machine Learning*. Citeseer.
 - [22] H Jin Kim, Michael I Jordan, Shankar Sastry, and Andrew Y Ng. 2004. Autonomous helicopter flight via reinforcement learning. In *Advances in neural information processing systems*. 799–806.
 - [23] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
 - [24] Jonathan Ko, Daniel J Klein, Dieter Fox, and Dirk Haehnel. 2007. Gaussian processes and reinforcement learning for identification and control of an autonomous blimp. In *Proceedings 2007 IEEE international conference on robotics and automation*. IEEE, 742–747.
 - [25] Petar Kormushev, Sylvain Calinon, and Darwin G Caldwell. 2013. Reinforcement learning in robotics: Applications and real-world challenges. *Robotics* 2, 3 (2013), 122–148.
 - [26] Hoang Le, Andrew Kang, Yisong Yue, and Peter Carr. 2016. Smooth imitation learning for online sequence prediction. In *International Conference on Machine Learning*. PMLR, 680–688.
 - [27] Sergey Levine, Zoran Popovic, and Vladlen Koltun. 2011. Nonlinear inverse reinforcement learning with gaussian processes. In *Advances in neural information processing systems*. 19–27.
 - [28] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971* (2015).
 - [29] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence* 41, 8 (2018), 1979–1993.
 - [30] Andrew Y Ng, Stuart J Russell, et al. 2000. Algorithms for inverse reinforcement learning. In *Icml*. 663–670.
 - [31] Matteo Papini, Andrea Battistello, and Marcello Restelli. [n.d.]. Safe Exploration in Gaussian Policy Gradient. ([n. d.]).
 - [32] Jason Pavis and Ronald Parr. 2011. Non-Parametric Approximate Linear Programming for MDPs. In *AAAI*.
 - [33] Martin L Puterman. 2014. Markov decision processes: discrete stochastic dynamic programming.
 - [34] Benjamin Recht. 2019. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems* 2 (2019), 253–279.
 - [35] Stéphane Ross and Drew Bagnell. 2010. Efficient reductions for imitation learning. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 661–668.
 - [36] Stéphane Ross and J Andrew Bagnell. 2014. Reinforcement and imitation learning via interactive no-regret learning. *arXiv preprint arXiv:1406.5979* (2014).
 - [37] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. 2011. 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*. 627–635.
 - [38] Stuart Russell. 1998. Learning agents for uncertain environments. In *Proceedings of the eleventh annual conference on Computational learning theory*. ACM, 101–103.
 - [39] Saumya Kumar Saksena, B Navaneethkrishnan, Sinchana Hegde, Pragadeesh Raja, and Ravi M Vishwanath. 2019. Towards Behavioural Cloning for Autonomous Driving. In *2019 Third IEEE International Conference on Robotic Computing (IRC)*. IEEE, 560–567.
 - [40] Stefan Schaal. 1999. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences* 3, 6 (1999), 233–242.
 - [41] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*. 1889–1897.
 - [42] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438* (2015).
 - [43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
 - [44] Qianli Shen, Yan Li, Haoming Jiang, Zhaoran Wang, and Tuo Zhao. 2020. Deep Reinforcement Learning with Smooth Policy. *arXiv preprint arXiv:2003.09534* (2020).
 - [45] Richard S Sutton, Andrew G Barto, et al. 1998. *Introduction to reinforcement learning*. Vol. 135. MIT press Cambridge.
 - [46] Andrew R Teel. 1995. Examples of stabilization using saturation: an input-output approach. *IFAC Proceedings Volumes* 28, 14 (1995), 209–214.
 - [47] Emanuel Todorov, Tom Erez, and Yuval Tassa. 2012. MuJoCo: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 5026–5033.
 - [48] Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. 2019. Unsupervised data augmentation for consistency training. *arXiv preprint arXiv:1904.12848* (2019).
 - [49] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. 2019. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573* (2019).
 - [50] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. 2008. Maximum Entropy Inverse Reinforcement Learning. In *AAAI*, Vol. 8. Chicago, IL, USA, 1433–1438.

A PROOFS

A.1 Proof of Theorem 1

PROOF. This theorem is a generalization of Theorem 5.9 in [32] for the case of continuous state and action space. The proof follows along the lines of one in [32] by replacing discrete Bellman optimality operator for \mathcal{V} and Q with continuous counterparts. Let $\mathcal{T}: \mathcal{V} \rightarrow \mathcal{V}$ be the Bellman optimality operator, where \mathcal{V} is the space of value functions. Then,

$$(\mathcal{T}V)(s) = \min_a \left(c(s, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s, a) V(s') ds' \right) \quad (19)$$

If $V(s)$ is L_v Lipschitz in s , then, so is $(\mathcal{T}V)(s)$. This can be seen as follows:

$\forall s_1, s_2 \in \mathcal{S}$ we have,

$$\begin{aligned} & |(\mathcal{T}V)(s_1) - (\mathcal{T}V)(s_2)| = \\ & \left| \min_a \left(c(s_1, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s_1, a) V(s') ds' \right) \right. \\ & \left. - \min_a \left(c(s_2, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s_2, a) V(s') ds' \right) \right|. \quad (20) \end{aligned}$$

Define \hat{a} as,

$$\hat{a} \triangleq \begin{cases} \arg \min_{a \in \mathcal{A}} \left(c(s_1, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s_1, a) V(s') ds' \right), \\ \quad \text{if } V(s_1) \leq V(s_2) \\ \arg \min_{a \in \mathcal{A}} \left(c(s_2, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s_2, a) V(s') ds' \right), \\ \quad \text{if } V(s_1) > V(s_2). \end{cases} \quad (21)$$

Then,

$$\begin{aligned} & |(\mathcal{T}V)(s_1) - (\mathcal{T}V)(s_2)| \leq \\ & \left| c(s_1, \hat{a}) + \gamma \int_{s' \in \mathcal{S}} p(s'|s_1, \hat{a}) V(s') ds' \right. \\ & \left. - c(s_2, \hat{a}) - \gamma \int_{s' \in \mathcal{S}} p(s'|s_2, \hat{a}) V(s') ds' \right|, \quad (22) \end{aligned}$$

where we have used the fact that $\forall f_1, f_2: \mathcal{S} \rightarrow \mathbb{R}$ bounded functions such that $\min_{s \in \mathcal{S}} f_1(s) \leq \min_{s \in \mathcal{S}} f_2(s)$ and $\hat{s} = \arg \min_{s \in \mathcal{S}} f_1(s)$, we have $|\min_{s \in \mathcal{S}} f_1(s) - \min_{s \in \mathcal{S}} f_2(s)| \leq |f_1(\hat{s}) - f_2(\hat{s})|$. Then using triangle inequality, we have $\forall s_1, s_2 \in \mathcal{S}$,

$$\begin{aligned} & |(\mathcal{T}V)(s_1) - (\mathcal{T}V)(s_2)| \leq \|c(s_1, \hat{a}) - c(s_2, \hat{a})\| + \\ & \left| \gamma \int_{s' \in \mathcal{S}} (p(s'|s_1, \hat{a}) - p(s'|s_2, \hat{a})) V(s') ds' \right|. \quad (23) \end{aligned}$$

$(1/L_v)V$ is 1-Lipschitz, hence using the Definition 3 and Lipschitz-ness of cost we get,

$$|(\mathcal{T}V)(s_1) - (\mathcal{T}V)(s_2)| \leq (L_c + L_v L_p) \|s_1 - s_2\|. \quad (24)$$

Hence, $(\mathcal{T}V)(s)$ is $L_c + \gamma L_v L_p$ -Lipschitz in s and a , if V is L_v -Lipschitz in s and a . Choose $V_0 = 0$. V^* is the stationary point of \mathcal{T} [33]: $V^* = \lim_{n \rightarrow \infty} \mathcal{T}^n V_0$. Hence, V^* is $L_c + \gamma L_p L_c + \gamma^2 L_p^2 L_c + \dots = L_c / (1 - \gamma L_p)$ -Lipschitz continuous.

Proof. b) we know that

$$Q^*(s, a) = c(s, a) + \gamma \int_{s' \in \mathcal{S}} p(s'|s, a) V^*(s') ds' \quad (25)$$

For this optimal state action value function Q^* and $s_1, s_2 \in \mathcal{S}$, we then have

$$\begin{aligned} & |Q^*(s_1, a_1) - Q^*(s_2, a_2)| \leq |c(s_1, a_1) - c(s_2, a_2)| + \\ & \left| \gamma \int_{s' \in \mathcal{S}} (p(s'|s_1, a_1) - p(s'|s_2, a_2)) V(s') ds' \right| \quad (26) \end{aligned}$$

$$\implies |Q^*(s_1, a_1) - Q^*(s_2, a_2)| \leq (L_c + \gamma L_v L_p) (\|s_1 - s_2\| + \|a_1 - a_2\|) \quad (27)$$

Using the value of L_v from the part a), we get

$$|Q^*(s_1, a_1) - Q^*(s_2, a_2)| \leq \left(\frac{L_c}{1 - \gamma L_p} \right) (\|s_1 - s_2\| + \|a_1 - a_2\|) \quad (28)$$

Hence, Proved. \square

A.2 Proof of Theorem 2

PROOF. From Lipschitz continuity of optimal Q^* function, we have $\forall s_1, s_2 \in \mathcal{S}$,

$$|Q^*(s_1, a_1) - Q^*(s_2, a_2)| \leq \left(\frac{L_c}{1 - \gamma L_p} \right) (\|s_1 - s_2\| + \|a_1 - a_2\|). \quad (29)$$

This can be re-written as

$$\left| \max_{a_1} Q^*(s_1, a_1) - \max_{a_2} Q^*(s_2, a_2) \right| \quad (30)$$

$$\leq \left(\frac{L_c}{1 - \gamma L_p} \right) (\|s_1 - s_2\| + \|a_1 - a_2\|). \quad (31)$$

Now, it is easy to see that $|\max_{a_1} Q^*(s_1, a_1) - \max_{a_2} Q^*(s_2, a_2)| \leq \max_a |Q^*(s_1, a) - Q^*(s_2, a)|$ [18]. From equivalence of norms, there exists a $\frac{1}{K_1} \geq 0$ such that

$$\max_a |Q^*(s_1, a) - Q^*(s_2, a)| \leq K_1 \left| \max_{a_1} Q^*(s_1, a_1) - \max_{a_2} Q^*(s_2, a_2) \right|.$$

Note that $\max_a |Q^*(s_1, a) - Q^*(s_2, a)|$ and κ are both pseudo-metrics on the space of $Q(s, \cdot)$ functions. Thus, from topological equivalence of pseudo-metrics, there exists a $K_2 \geq 0$ such that $\kappa(Q(s_1, \cdot), Q(s_2, \cdot)) \leq K_2 \max_a |Q(s_1, a) - Q(s_2, a)|$. Thus, we have $\kappa(Q(s_1, \cdot), Q(s_2, \cdot)) \leq K_1 K_2 \left| \max_{a_1} Q^*(s_1, a_1) - \max_{a_2} Q^*(s_2, a_2) \right|$. Thus, for all $s_1, s_2 \in \mathcal{S}$ and for all $a \in \mathcal{A}$ we have,

$$\kappa(Q(s_1, \cdot), Q(s_2, \cdot)) \leq K_1 K_2 \left(\frac{L_c}{1 - \gamma L_p} \right) (\|s_1 - s_2\|). \quad (32)$$

Now, using definition and Lipschitz continuity of H , we get

$$\|\mu(s_1) - \mu(s_2)\| = \|H(Q^*(s_1, \cdot)) - H(Q^*(s_2, \cdot))\| \quad (33)$$

$$\leq L_\mu \kappa(Q^*(s_1, \cdot), Q^*(s_2, \cdot)) \quad (34)$$

$$\leq K_1 K_2 L_\mu \frac{L_c}{1 - \gamma L_p} \|s_1 - s_2\|. \quad (35)$$

Thus, the optimal mean policy μ^* is $\frac{K_1 K_2 L_\mu L_c}{1 - \gamma L_p}$ -Lipschitz continuous with respect to the states. Then, the stationary stochastic policy π^* obtained as $\mathcal{N}(\mu^*(s), \sigma)$ (for a fixed σ) is $\frac{K_1 K_2 L_\mu L_c}{1 - \gamma L_p}$ -Lipschitz continuous with respect to the states. \square

B EXPERIMENTAL SETUP AND HYPERPARAMETER DETAILS

B.1 Policy, Cost and Value Models

Policy. In our work, the policy, $\pi_\theta(a | s)$, is parameterized with a neural network. The neural network takes in a state vector s and deterministically maps s to a vector μ . The neural network is additionally used to learn a vector r of log standard deviation with the same dimension as a . During training, the action is sampled stochastically from $\mathcal{N}(\mu, \exp(r))$. While evaluating a policy, the mean action μ is chosen.

Value. The value neural network simply takes in a state s and outputs scalar $V(s)$.

Discriminator. The discriminator network takes in a state-action pair (s, a) and outputs a scalar x . The cost of this state-action pair is then defined as $c(s, a) = \log(S(x))$, where $S(x) = \frac{e^x}{e^x + 1}$.

B.1.1 Neural network details. All the models (policy, value, and cost functions) in this work are multi-layer perceptrons (MLPs). The policy models for all the algorithms (TRPO, GAIL, and SPaCIL) have two layers of 400 and 300 neurons each. The value functions across all the algorithms have two layers of 100 neurons each. Similarly, the discriminator network (a representative for cost model) has two layers of 100 neurons each. The non-linearity used in all the models is *tanh*. All the weights and biases are initialized using $U(-k, k)$ where $k = \frac{1}{\sqrt{\text{size of weights}}}$. The final layer across models has zero bias and weights multiplied by 0.1. The step size and learning rate for the policy model is determined by the TRPO algorithm. All the other models are trained using Adam optimizer [23] with a learning rate of $1e-3$ for the value network, and a learning rate of 0.01 for the discriminator network.

B.2 Expert learning using TRPO

The MuJoCo environments used in our experiments are from the OpenAI Gym [6]. We specifically work with v2 of Reacher, Hopper, Walker2d, HalfCheetah, and Ant. The expert is trained using TRPO [41]. The TRPO hyperparameters *max KL* and *damping* are both fixed to 0.01 for all the environments. The batch size is 50000 for all the environments except Reacher’s 5000. We train all the environments for 500 iterations except Reacher’s 200. We fix $\gamma = 0.995$ for all the environments except Reacher’s 0.99. We fix $\tau = 0.95$ for Reacher, $\tau = 0.99$ for Hopper and Walker, and $\tau = 0.97$ for the rest.

B.3 GAIL and SPaCIL parameters

The best set of hyperparameters for SPaCIL are listed in Table 3. The policy regularization strength is denoted by λ_1 . The cost regularization strength is denoted by λ_2 . The step size parameter in the projected gradient descent part of regularizer estimation is denoted as ν and takes a value of 0.02 across environments. The perturbation strength around a particular state (to project onto ϵ -Ball) is denoted by ϵ . It takes a value of 0.01 across environments. For regularization purposes, we keep the policy σ fixed to 1. The mixing parameter μ is randomly sampled from a uniform distribution over $[0, 1)$. sampled Generalized advantage estimation (GAE, [42]) parameters, γ and τ are same as TRPO for all the environments.

GAIL has all the hyperparameters, except training data size, same as that of TRPO. We are provided with a fixed number of expert trajectories before training begins. This number is the same for both GAIL and SPaCIL, and is listed in Table 3. Each trajectory consists of 50 (s, a) pairs for Reacher-v2, and 1000 (s, a) pairs for all the other tasks.

Table 3: Environments details and performance of expert policies.

Environment	λ_1	λ_2	Expert traj No.	Agent traj No.
Reacher-v2	0.01	0.001	50	50
Hopper-v2	0.001	0.001	6	6
Walker2d-v2	0.001	0.001	10	10
HalfCheetah-v2	0.01	0.001	6	6
Ant-v2	0.001	0.001	15	15

B.4 The MuJoCo Tasks

Reacher-v2. In this environment, the agent is a grasping arm with a hinge, a body with a joint, and a tip that grasps (see Fig. 5). With the hinge fixed at the center of a square grid environment, the agent’s tip is placed at a random starting state. The agent’s goal is to be able to grasp a target that is randomly spawned in the square grid.

Hopper-v2. In this environment, the agent is a robot with a torso and a leg (see Fig. 5). The agent is tasked with learning to hop through the environment. The learned walking behaviour is desired to have a stable gait.

HalfCheetah-v2. In this environment, the agent is a robot with only one forelimb and one hind-limb (see Fig. 5). The agent is tasked with learning to walk and hop through the environment. The learned walking behaviour is desired to have a stable gait.

Walker2d-v2. In this environment, the agent is a robot with a torso and two legs (see Fig. 5). The agent is tasked with learning to jump and walk through the environment. The learned walking behaviour is desired to have a stable gait.

Ant-v2. In this environment, the agent is a robot with four legs (see Fig. 5). The agent is tasked with learning to hop through the environment. The learned walking behaviour is desired to have a stable gait.

B.5 Deep RL tricks used in the implementation

We use the following standard tricks from deep reinforcement learning (DRL) in our implementation:

- (1) We keep and use a running average of the states to deal with covariate shift in the input data.
- (2) We evaluate the performance of our algorithms on a separate test environment. We keep track of the best evaluation return to save our best agent model. We sample 20,000 s - a pairs at each test iteration for testing. The action is sampled as the mean action ($\mu(s)$) at a certain state (s) rather than from the stochastic policy. The evaluation performance curves for GAIL and SPaCIL are included in Fig. 6).

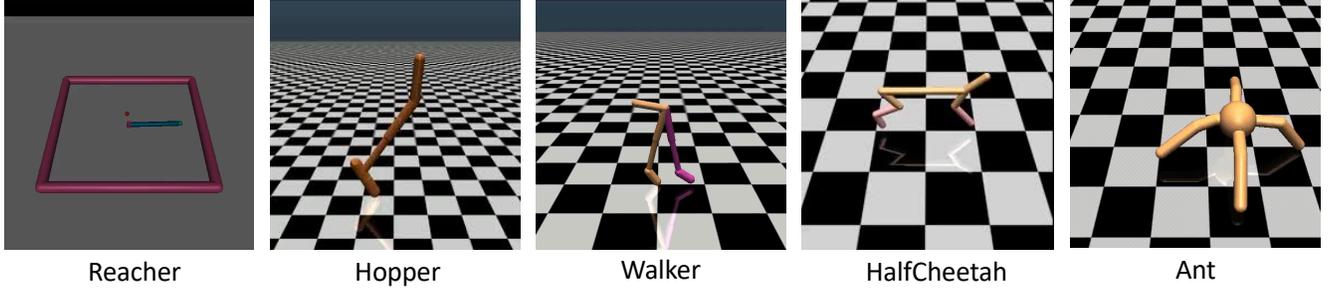


Figure 5: Robotic continuous control tasks from MuJoCo.

- (3) Different components of our code are run on different devices (*i.e.*, CPU and GPU) to optimize for algorithm’s run time. While sampling trajectories from the agent policy, we use multi-threading and run the code on the CPU. All the gradient calculations vis back-propagation are performed on GPU cores.

C IMPLEMENTATION OF PROJECTION

At the crux of our algorithm are the following policy and cost regularizers:

$$\mathcal{R}_s^\pi(\theta_k) = \mathbb{E}_{s \sim \rho^{\pi_{\theta_k}}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{D}_J(\pi_{\theta_k}(s), \pi_{\theta_k}(\tilde{s})), \quad (36)$$

and

$$\mathcal{R}_{(s,a)}^c(\theta) = \mathbb{E}_{s \sim \rho^{\pi_\theta}} \max_{\tilde{s} \in \mathbb{B}_d(s, \epsilon)} \|c(s, \pi_\theta(\cdot|s)) - c(s, \pi_\theta(\cdot|\tilde{s}))\|_2. \quad (37)$$

Both the regularizers in Eqn. 36 and Eqn. 37 require us to solve a maximization problem over the ϵ -ball around a certain state s . The goal of this maximization problem is to find a state s' within an ϵ -ball of a state s at which a certain function, $f(s, s')$ takes the maximum value. In Algorithm 3, we discuss a general projected gradient based approach to solve this maximization problem that is applicable to both Eqn. 36 and Eqn. 37. For the policy smoothing regularizer of Eqn. 36 the f in Algorithm 3 is given by $f(s, s') = \mathcal{D}_J(\pi_\theta(s), \pi_\theta(s'))$ (the Jeffrey’s divergence between policies $\pi_\theta(s)$ and $\pi_\theta(s')$). For the cost smoothing regularizer of Eqn. 14, $f(s, s') = \|c(s, \pi_\theta(\cdot|s)) - c(s, \pi_\theta(\cdot|s'))\|_2$ (the L_2 distance between costs). Here we provide exact details of how to obtain $\nabla_\delta f(s, s + \delta_\ell)$ (in step 6) and how to project onto the ball $\Pi_{\mathbb{B}_d}$ (step 8).

For $f(s, s') = \mathcal{D}_J(\pi_\theta(s), \pi_\theta(s'))$ and Gaussian distributed policies (*i.e.*, $\pi_\theta(\cdot|s) \stackrel{d}{=} \mathcal{N}(\mu_\theta(s), \sigma^2)$), f is reduced to $\frac{\|\mu_\theta(s) - \mu_\theta(s')\|_2^2}{\sigma^2}$. Gradient of this quantity is estimated using automatic back-propagation [15] through the policy neural network. Gradient of $f(s, s') = \|c(s, \pi_\theta(\cdot|s)) - c(s, \pi_\theta(\cdot|s'))\|_2$ can be equivalently estimated using back-propagation.

Once we have $\delta_{\ell+1}$ from Step 6 of Algorithm 3, the projection onto the ball $\mathbb{B}_d(0, \epsilon)$ is evaluated using the following formula:

$$\delta_{\text{new}} = \delta_{\text{old}} \min\left\{1, \frac{\epsilon}{\|\delta_{\text{old}}\|_2}\right\}. \quad (38)$$

Algorithm 3 Maximization of f over ϵ -ball of a state s

- 1: **Input:** s, ϵ, η_δ
 - 2: **Initialize:** δ_0
 - 3: (N steps of projected gradient descent)
 - 4: **for** $\ell = 1, 2, \dots, N-1$
 - 5: **(Update δ in the direction of increase in f)**
 - 6: $\delta_{\ell+1} = \delta_\ell + \eta_\delta \nabla_\delta f(s, s + \delta_\ell)$
 - 7: **(Project $\delta_{\ell+1}$ onto the ϵ -ball)**
 - 8: $\delta_{\ell+1} = \Pi_{\mathbb{B}_d(0, \epsilon)}(\delta_{\ell+1})$
 - 9: **end for**
-

D DISCUSSION ON SMOOTHNESS METRIC

An alternative view of smoothness metric. To quantify the smoothness of the learned policy, we had introduced a novel metric that took the following form:

$$J(\{\tau_i\}_{i=1}^N) = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^T \max_{\tilde{s}_j \in \mathbb{B}_d(s_j, \epsilon)} \frac{\|\mu(s_j) - \mu(\tilde{s}_j)\|}{\|s - \tilde{s}_j\|}, s_j \sim \tau_i. \quad (39)$$

where $\mu(s)$ is the deterministic mean function for $\pi(\cdot|s)$. We had gotten the above form for the metric by starting out by defining Lipschitz constant L corresponding to the deterministic mean function, $\mu(s)$ for a general norm is given by

$$L = \max_{s_1 \neq s_2} \frac{\|\mu(s_1) - \mu(s_2)\|}{\|s_1 - s_2\|}, \quad s_1, s_2 \in \mathcal{S}. \quad (40)$$

Taking $s_1 = s + \delta s$ and $s_2 = s$

$$L = \max_{\delta s \neq 0} \frac{\|\mu(s + \delta s) - \mu(s)\|}{\|s + \delta s - s\|}, \quad (41)$$

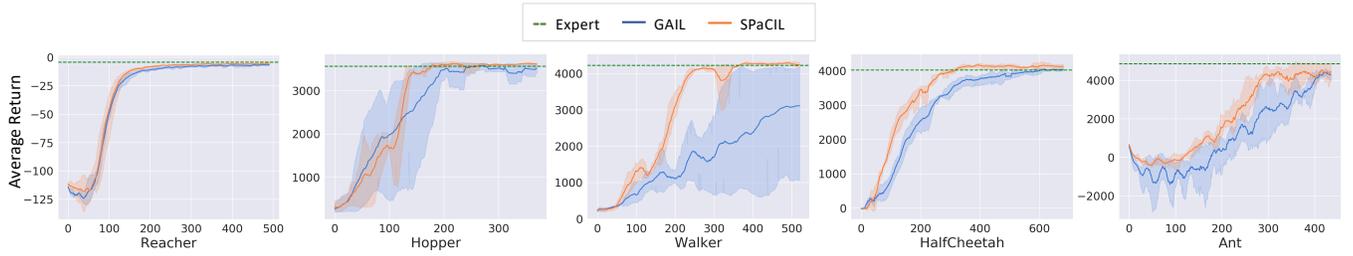


Figure 6: Evaluation curves for GAIL and SPaCIL on challenging continuous control tasks. The dotted green line represents the average return of the demonstration dataset. The dual regularization of policies and costs results in faster learning for SPaCIL.

Table 4: Validating our smoothness metric is a good measure of smoothness

	0	0.1	0.09	0.05	0.01	0.001
Reacher G	-4.18±1.79	-163.44±43.95	-47.07±11.08	-30.12±10.25	-11.67±2.18	-12.35±2.14
Reacher J	2.69e-5±2.02e-5	14.42±2.06	1.02±0.12	0.29±0.05	8.67e-4±1.79e-4	3.72e-4±7.85e-5
Hopper G	3662.71±39.27	79.73±73.48	4.69±0.033	383.59±313.75	2644.44±1207.97	3667.35±9.72
Hopper J	7.77±0.067	18.24±3.01	18.34±0.0162	5.43±0.13	6.732±0.27	7.86±0.19
Walker2d G	4397.58±6.92	3.21±2.81	61.92±53.04	245.21±119.8	4199.56±17.19	4329.44±10.69
Walker2d J	35.61±0.067	85.91±1.06	125.97±6.06	37.56±1.60	39.37±1.1	38.53±0.12
HalfCheetah G	4141.10±93.84	-743.56±60.59	-729.94±242.07	128.97±141.52	3642.23±75.71	4061.67±75.59
HalfCheetah J	17.73±0.11	61.63±6.72	88.67±33.76	4.47±2.03	17.72±0.19	18.51±0.21
Ant G	4788.73±75.18	-4880.08±1719.81	-1454.50±961.96	1685.58±482.09	4561.18±617.04	4628.54±893.13
Ant J	0.93±0.09	186.91±74.77	43.58±27.72	1.64±2.27	0.91±0.07	0.97±0.04

Using the first order Taylor series approximation for μ at s : $\mu(s + \delta s) = \mu(s) + J_{\mu(s)}\delta s$ we get

$$L = \max_{\delta s_2 \neq 0} \frac{\|J_{\mu(s)}\delta s\|}{\|\delta s\|}, \quad (42)$$

For L^2 norm, the quantity on the right in Eqn. 42 is the spectral norm of $J_{\mu(s)}$. Hence, the local Lipschitz constant, L at a particular state s is given by the spectral norm of the Jacobian, J at that state: $\|J_{\mu(s)}\|_2 = \sigma_{\max}(J_{\mu(s)})$, *i.e.*, the maximum singular value [10, 46]. Therefore, another approach to evaluate the smoothness of π , is to estimate the expected Jacobian norm using sampled trajectories as,

$$\hat{E}[\|J_{\mu(s)}\|_2] = \frac{1}{NT} \sum_{i=1}^N \sum_{j=1}^T \|J_{\mu(s_j)}\|_2, \quad s_j \sim \tau_i \quad (43)$$

where T and N are the number of sampled trajectories and time steps, respectively. A sampled trajectory, $\tau \sim \pi$ is a trajectory of the form $\{s_0, a_0, s_1, a_1, \dots, s_T\}$, where $s_0 \sim \rho_0$ is the starting state, $a_t \sim \pi(\cdot|s_t)$, and T denotes the time step at which we terminate an episode. The quantity in Eqn. 43 is the average spectral norm instead of the maximum of the norms. Hence, in our work, we use Eqn. 39 as the smoothness metric. $J(\{\tau_i\}_{i=1}^N)$ in Eqn. 39 can be estimated using the batch $\{\tau_i\}_{i=1}^N$ of data sampled from any policy.

E MORE RESULTS

We include results from some more experiments here. Fig. 6 shows the average return (and one standard deviation) over an evaluation environment. Table 4 reports average return and average metric for perturbed models various tasks where perturbation is the variance of Gaussian noise added to the original model parameters.