
TaSIL: Taylor Series Imitation Learning

Daniel Pfrommer*

Massachusetts Institute of Technology[†]
Cambridge, MA
dpfrom@mit.edu

Thomas T.C.K. Zhang*

University of Pennsylvania
Philadelphia, PA
ttz2@seas.upenn.edu

Stephen Tu

Robotics at Google
New York, NY
stephentu@google.com

Nikolai Matni

University of Pennsylvania
Philadelphia, PA
nmatni@seas.upenn.edu

Abstract

We propose Taylor Series Imitation Learning (TaSIL), a simple augmentation to standard behavior cloning losses in the context of continuous control. TaSIL penalizes deviations in the higher-order Taylor series terms between the learned and expert policies. We show that experts satisfying a notion of *incremental input-to-state stability* are easy to learn, in the sense that a small TaSIL-augmented imitation loss over expert trajectories guarantees a small imitation loss over trajectories generated by the learned policy. We provide generalization bounds for TaSIL that scale as $\tilde{O}(1/n)$ in the realizable setting, for n the number of expert demonstrations. Finally, we demonstrate experimentally the relationship between the robustness of the expert policy and the order of Taylor expansion required in TaSIL, and compare standard Behavior Cloning, DART, and DAgger with TaSIL-loss-augmented variants. In all cases, we show significant improvement over baselines across a variety of MuJoCo tasks.

1 Introduction

Imitation learning (IL), wherein expert demonstrations are used to train a policy [1, 2], has been successfully applied to a wide range of tasks, including self-driving cars [3, 4], robotics [5], and video game playing [6]. While IL is typically more sample-efficient than reinforcement learning-based alternatives, it is also known to be sensitive to distribution shift: small errors in the learned policy can lead to compounding errors and ultimately, system failure [3, 6]. In order to mitigate the effects of distribution shift caused by policy error, two broad approaches have been taken by the community. On-policy approaches such as DAgger [6] augment the data-set with expert-labeled or corrected trajectories generated by learned policies. In contrast, off-policy approaches such as DART [7] and GAIL [8] augment the data-set by perturbing the expert controlled system. In both cases, the goal is to provide examples to the learned policy of how the expert recovers from errors. While effective, these methods either require an interactive expert (on-policy), or access to a simulator for policy rollouts during training (off-policy), which may not always be practically feasible.

In this work, we take a more direct approach towards mitigating distribution shift. Rather than providing examples of how the expert policy recovers from error, we seek to endow a learned policy with the robustness properties of the expert directly. In particular, we make explicit the underlying assumption in previous work that a good expert is able to recover from perturbations through

*Both authors contributed equally to this work.

[†]This work was done while affiliated with the University of Pennsylvania.

the notion of *incremental input-to-state stability*, a well-studied control theoretic notion of robust nonlinear stability. Under this assumption, we show that if the p -th order Taylor series approximation of the learned and expert policies approximately match on expert-generated trajectories, then the learned policy induces closed-loop behavior similar to that of the expert closed-loop system. Here the order p is determined by the robustness properties of the expert, and makes quantitative the informal observation that more robust experts are easier to learn. Fundamentally, we seek to characterize settings where offline imitation learning is Probably Approximately Correctly (PAC)-learnable [9]; that is, defining notions of expert data and designing algorithms such that low training time error on offline expert demonstrations implies low test time error along the learned policy’s trajectories, in spite of distribution shift.

Contributions We propose and analyze Taylor Series Imitation Learning (TaSIL), which augments the standard behavior cloning imitation loss to capture errors between the higher-order terms of the Taylor series expansion of the learned and expert policies. We reduce the analysis of TaSIL to the analysis of a supervised learning problem over expert data: in particular, we identify a robustness criterion for the expert such that a small TaSIL-augmented imitation loss over expert trajectories guarantees that the difference between trajectories generated by the learned and expert policies is also small. We also provide a finite-difference based approximation that is applicable to experts that cannot directly query their higher-order derivatives, expanding the practical applicability of TaSIL. We show in the realizable setting that our algorithm achieves generalization bounds scaling as $\tilde{O}(1/n)$, for n the number of expert demonstrations. Finally, we empirically demonstrate (i) that the relationship between the robustness of the expert policy and the order of Taylor expansion required in TaSIL predicted by our theory is observed in practice, and (ii) the benefits of using the TaSIL-augmented imitation loss by comparing the sample-efficiency of standard and TaSIL-loss-augmented behavior cloning, DART, and DAgger on a variety of MuJoCo tasks: on hard instances where behavior cloning fails, the TaSIL-augmented variants show significant performance gains.

1.1 Related work

Imitation learning Behavior cloning is known to be sensitive to compounding errors induced by small mismatches between the learned and expert policies [3, 6]. On-policy [6] and off-policy [7, 8] approaches exist that seek to prevent this distribution shift by augmenting the data-set created by the expert. While DAgger is known to enjoy $\tilde{O}(T)$ sample-complexity in the task horizon T for loss functions that are strongly convex in the policy parameters,³ we are not aware of finite-data guarantees for DART or GAIL. In addition to these seminal papers, there is a body of work that seeks to leverage control theoretic techniques [10, 11] to ensure (robust) stability of the learned policy. More closely related to our work are the results by Ren et al. [12] and Tu et al. [13]. In Ren et al. [12], a two-stage pipeline of imitation learning followed by policy optimization via a PAC-Bayes regularization term is used to provide generalization bounds of the learned policy across randomly drawn environments. This work is mostly complementary to ours, as TaSIL could in principle be used to augment the imitation losses used in the first stage of their pipeline (leaving the second stage unmodified).

In Tu et al. [13], sample-complexity bounds for IL are provided under the assumption that the learned policy can be explicitly constrained to match the incremental stability properties of the expert. While conceptually appealing, practically enforcing such stability constraints is difficult, and as such Tu et al. [13] resort to heuristics in their implementation. In contrast, we provide sample-complexity guarantees for a practically implementable algorithm, and under much milder stability assumptions.

Robust stability and learning for continuous control There is a rich body of work applying Lyapunov stability or contraction theory [14] to learning for continuous control. For example [15–18] use stability-based regularizers to trim the hypothesis space, and empirically show that this leads to more sample-efficient and robust learning methods. Lyapunov stability and contraction theory have also been used to provide finite sample-complexity guarantees for adaptive nonlinear control [19] and learning stability certificates [20].

³This bound degrades to $\tilde{O}(T^2)$ when loss function is only convex, and does not hold for the nonconvex loss functions we consider.

2 Problem formulation

We consider the nonlinear discrete-time dynamical system

$$x_{t+1} = f(x_t, u_t), x_0 = \xi, \quad (1)$$

where $x_t \in \mathbb{R}^d$ is the system state, $u_t \in \mathbb{R}^m$ is the control input, and $f : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$ defines the system dynamics. We study system (1) evolving under the control input $u_t = \pi(x_t) + \Delta_t$, for $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ a suitable control policy, and $\Delta_t \in \mathbb{R}^m$ an additive input perturbation (that will be used to capture policy errors). We define the corresponding *perturbed* closed-loop dynamics by $x_{t+1} = f_{\text{cl}}^{\pi}(x, \Delta_t) := f(x, \pi(x) + \Delta_t)$, and use $x_t^{\pi}(\xi, \{\Delta_s\}_{s=0}^{t-1})$ to denote the value of the state x_t at time t evolving under control input $u_t = \pi(x) + \Delta_t$ starting from initial condition $x_0 = \xi$. To lighten notation, we use $x_t^{\pi}(\xi)$ to denote $x_t^{\pi}(\xi, \{0\}_{s=0}^{t-1})$, and overload $\|\cdot\|$ to denote the Euclidean norm for vectors and the operator norm for matrices and tensors.

Initial conditions ξ are assumed to be sampled randomly from a distribution \mathcal{D} with support restricted to a compact set \mathcal{X} . We assume access to n rollouts of length T from an expert π_* , generated by drawing initial conditions $\{\xi_i\}_{i=1}^n$ i.i.d. from \mathcal{D} . The IL task is to learn a policy $\hat{\pi}$ which leads to a closed-loop system with similar behavior to that induced by the expert policy π_* as measured by the *expected imitation gap* $\mathbb{E}_{\xi} \max_{1 \leq t \leq T} \|x_t^{\hat{\pi}}(\xi) - x_t^{\pi_*}(\xi)\|$.

The baseline approach to IL, typically referred to as *behavior cloning* (BC), casts the problem as an instance of supervised learning. Denote the discrepancy between an evaluation policy $\bar{\pi}$ and the expert policy π_* on a trajectory generated by a rollout policy π_d starting at initial condition ξ by $\Delta_t^{\pi_d}(\xi; \bar{\pi}) := \bar{\pi}(x_t^{\pi_d}(\xi)) - \pi_*(x_t^{\pi_d}(\xi))$. BC directly solves the supervised empirical risk minimization (ERM) problem

$$\hat{\pi}_{\text{bc}} \in \operatorname{argmin}_{\pi \in \Pi} \frac{1}{n} \sum_{i=1}^n h(\{\Delta_t^{\pi_*}(\xi_i; \pi)\}_{t=0}^{T-1}),$$

over a suitable policy class Π . Here, $h : (\mathbb{R}^m)^T \rightarrow \mathbb{R}$ is a loss function which encourages the discrepancy terms $\Delta_t^{\pi_*}(\xi_i; \pi)$ to be small along the expert trajectories. While behavior cloning is conceptually simple, it can perform poorly in practice due to distribution shifts triggered by errors in the learned policy $\hat{\pi}_{\text{bc}}$. Specifically, due to the effects of compounding errors, the closed-loop system induced by the behavior cloning policy $\hat{\pi}_{\text{bc}}$ may lead to a dramatically different distribution over system trajectories than that induced by the expert policy π_* , even when the population risk on the expert data, $\mathbb{E}_{\xi} h(\{\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{bc}})\}_{t=0}^{T-1})$, is small.

As described in the introduction, existing approaches to mitigating distribution shift seek to augment the data with examples of the expert recovering from errors. These approaches either require an interactive oracle (e.g., DAgger) or access to a simulator for policy rollouts (e.g., DART, GAIL), and may not always be practically applicable. To address the distribution shift challenge without resorting to data-augmentation, we propose TaSIL, an off-policy IL algorithm which provably leads to learned policies that are robust to distribution shift.

The rest of the paper is organized as follows: in Section 3, we focus on ensuring robustness to policy errors for a single initial condition ξ . We show that the imitation gap $\|x_t^{\pi}(\xi) - x_t^{\pi_*}(\xi)\|$ between a test policy π and the expert policy π_* can be controlled by the TaSIL-augmented imitation loss evaluated on the expert trajectory $\{x_t^{\pi_*}(\xi)\}$, effectively reducing the analysis of the imitation gap to a supervised learning problem over the expert data. In Section 4, we integrate these results with tools from statistical learning theory to show that $n \gtrsim \varepsilon^{-r}/\delta$ trajectories are sufficient to achieve imitation gap of at most ε with probability at least $1 - \delta$; here $r > 0$ is a constant determined by the stability properties of the expert policy π_* , with more robust experts corresponding to smaller values of r . Finally, in Section 5, we validate our analysis empirically, and show that using the TaSIL-augmented loss function in IL algorithms leads to significant gains in performance and sample efficiency.

3 Bounding the imitation gap on a single trajectory

In this section, we fix an initial condition ξ and test policy π , and seek to control the imitation gap $\Gamma_T(\xi; \pi) := \max_{1 \leq t \leq T} \|x_t^{\pi}(\xi, \{0\}) - x_t^{\pi_*}(\xi, \{0\})\|$. A natural way to compare the

closed-loop behavior of a test policy π to that of the expert policy π_* is to view the discrepancy $\Delta_t^\pi(\xi; \pi) := \pi(x_t^\pi(\xi)) - \pi_*(x_t^\pi(\xi))$ as an *input perturbation* to the expert closed-loop system. By writing $x_{t+1} = f_{\text{cl}}^\pi(x_t, 0) = f_{\text{cl}}^{\pi_*}(x_t, \Delta_t^\pi(\xi; \pi))$, the imitation gap can be written as $\Gamma_T(\xi; \pi) = \max_t \|x_t^{\pi_*}(\xi, \{\Delta_s^\pi(\xi; \pi)\}_{s=0}^{t-1}) - x_t^{\pi_*}(\xi, \{0\})\|$, suggesting that closed-loop expert systems that are robust to input perturbations, as measured by the difference between nominal and perturbed trajectories, will lead to learned policies that enjoy smaller imitation gaps.

Stability conditions defined in terms of differences between nominal and perturbed trajectories have been extensively studied in robust nonlinear control theory via the notion of *incremental input-to-state stability* (δ -ISS) (see e.g., Angeli [21] and references therein). Before proceeding, we recall definitions of standard comparison functions [22]: a function $\gamma(x)$ is class \mathcal{K} if it is continuous, strictly increasing, and satisfies $\gamma(0) = 0$, and a function $\beta(x, t)$ is class \mathcal{KL} if it is continuous, $\beta(\cdot, t)$ is class \mathcal{K} for each fixed t , and $\beta(x, \cdot)$ is decreasing for each fixed x .

Definition 3.1 (δ -ISS system). *Consider the closed-loop system evolving under policy π and subject to perturbations Δ_t given by $x_{t+1} = f_{\text{cl}}^\pi(x_t, \Delta_t)$. The closed-loop system $f_{\text{cl}}^\pi(x_t, \Delta_t)$ is incremental input-to-state stable (δ -ISS) if there exists a class \mathcal{KL} function β and a class \mathcal{K} function γ such that for all initial conditions $\xi_1, \xi_2 \in \mathcal{X}$, perturbation sequences $\{\Delta_t\}_{t \geq 0}$, and $t \in \mathbb{N}$:*

$$\|x_t^\pi(\xi_1; \{\Delta_s\}_{s=0}^{t-1}) - x_t^\pi(\xi_2; \{0\}_{s=0}^{t-1})\| \leq \beta(\|\xi_1 - \xi_2\|, t) + \gamma\left(\max_{0 \leq k \leq t-1} \|\Delta_k\|\right). \quad (2)$$

Definition 3.1 says that: (i) trajectories generated by δ -ISS systems converge towards each other if they begin from different initial conditions, and (ii) the effect of bounded perturbations $\{\Delta_t\}$ on trajectories is bounded. Our results will only require the stability conditions of Definition 3.1 to hold for a class of norm bounded perturbations. In light of this, we say that a system is η -locally δ -ISS if equation (2) holds for all input perturbations satisfying $\sup_{t \in \mathbb{N}} \|\Delta_t\| \leq \eta$.

By writing $x_{t+1} = f_{\text{cl}}^\pi(x_t, 0) = f_{\text{cl}}^{\pi_*}(x_t, \Delta_t^\pi(\xi; \pi))$, we conclude that if $x_{t+1} = f_{\text{cl}}^{\pi_*}(x_t, \Delta_t)$ is η -locally δ -ISS, and that $\sup_{t \in \mathbb{N}} \|\Delta_t^\pi(\xi; \pi)\| \leq \eta$, then equation (2) yields

$$\|x_t^\pi(\xi) - x_t^{\pi_*}(\xi)\| \leq \gamma\left(\max_{0 \leq k \leq t-1} \|\Delta_k^\pi(\xi; \pi)\|\right) \implies \Gamma_T(\xi; \pi) \leq \gamma\left(\max_{0 \leq t \leq T-1} \|\Delta_t^\pi(\xi; \pi)\|\right). \quad (3)$$

Equation (3) shows that the imitation gap $\|x_t^\pi(\xi) - x_t^{\pi_*}(\xi)\|$ is controlled by the maximum discrepancy $\|\Delta_t^\pi(\xi; \pi)\|$ incurred on trajectories generated by the policy π . A natural way of bounding the test discrepancy $\|\Delta_t^\pi(\xi; \pi)\|$, defined over trajectories generated by π , by the training discrepancy $\|\Delta_t^{\pi_*}(\xi; \pi)\|$, defined over trajectories generated by π_* , is to write out a Taylor series expansion of the former around the latter, i.e., to write:⁴

$$\|\Delta_t^\pi(\xi; \pi)\| \leq \|\Delta_t^{\pi_*}(\xi; \pi)\| + \|\partial_x \Delta_t^{\pi_*}(\xi; \pi)\| \Gamma_T(\xi; \pi) + \mathcal{O}(\Gamma_T^2(\xi; \pi)), \quad (4)$$

where $\partial_x \Delta_t^{\pi_*}(\xi; \pi)$ denotes the partial derivative of the discrepancy with respect to the argument of the policies π and π_* . The challenge however, is that the imitation gap $\Gamma_T(\xi; \pi)$ also appears in the Taylor series expansion (4), leading to an implicit constraint. We show that this can be overcome if the order p of the Taylor series expansion (4) is sufficiently large, as determined by the decay rate of the class \mathcal{K} function $\gamma(\cdot)$ defining the robustness of the expert policy. We begin by identifying a condition reminiscent of adversarially robust training objectives that ensures a small imitation gap.

Proposition 3.1. *Let the expert closed-loop system $f_{\text{cl}}^{\pi_*}$ be η -locally δ -ISS for some $\eta > 0$. Fix an imitation gap bound $\varepsilon > 0$, initial condition ξ , and policy π . Then if*

$$\max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_*(x_t^{\pi_*}(\xi) + \delta) - \pi(x_t^{\pi_*}(\xi) + \delta)\| \leq \min\{\eta, \gamma^{-1}(\varepsilon)\}, \quad (5)$$

we have that the imitation gap satisfies $\Gamma_T(\xi; \pi) \leq \varepsilon$.

Proposition 3.1 states that if a policy π is sufficiently close to the expert policy π_* in a tube around the expert trajectory, then the imitation gap remains small. How to ensure that inequality (5) holds using only offline data from expert trajectories is not immediately obvious. The Taylor series expansion (4) suggests that a natural approach to satisfying this condition is to match derivatives of the test policy π to those of the expert policy π_* . We show next that a sufficient order for such a Taylor series

⁴We only take a first order expansion here for illustrative purposes.

expansion is naturally determined by the decay rate of the class \mathcal{K} function $\gamma(x)$ towards 0. Less robust experts have functions that decay to 0 more slowly and will lead to stricter sufficient conditions. Conversely, more robust experts have functions that decay to 0 more quickly, and will lead to more relaxed sufficient conditions. We focus on two disjoint classes of class \mathcal{K} functions: (i) functions that decay in their argument faster than a linear function, i.e. $\gamma(x) < \mathcal{O}(x)$ as $x \rightarrow 0^+$, and (ii) functions that decay in their argument no faster than a linear function, i.e. $\gamma(x) \geq \Omega(x)$ as $x \rightarrow 0^+$.

Rapidly decaying class \mathcal{K} functions

We show that when the class \mathcal{K} function $\gamma(x)$ decays to 0 faster than $\mathcal{O}(x)$ in some neighborhood of 0, then matching the *zeroth-order difference* $\max_t \|\Delta_t^{\pi_*}(\xi; \pi)\|$ on the expert trajectory, as is done in vanilla behavior cloning, suffices to close the imitation gap. We make the following assumption on the test policy π and expert policy π_* .

Assumption 3.1. *There exists a non-negative constant L_π such that $\|\bar{\pi}(x) - \bar{\pi}(y)\|_2 \leq L_\pi \|x - y\|_2$ for all x, y , and $\bar{\pi} \in \{\pi, \pi_*\}$.*

Proposition 3.1 then leads to the following guarantee on the imitation gap.

Theorem 3.1. *Fix a test policy π and initial condition $\xi \in \mathcal{X}$, and let Assumption 3.1 hold. Let $f_{\text{cl}}^{\pi_*}$ be η -locally δ -ISS for some $\eta > 0$, and assume that the class \mathcal{K} function $\gamma(\cdot)$ in (2) satisfies $\gamma(x) \leq \mathcal{O}(x^{1+r})$ for some $r > 0$. Choose constants $\mu, \alpha > 0$ such that*

$$2L_\pi x + (x/\mu)^{\frac{1}{1+r}} \leq \gamma^{-1}(x) \text{ for all } 0 \leq x \leq \alpha. \quad (6)$$

Provided that the imitation error on the expert trajectory incurred by π satisfies:

$$\max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi_*}(\xi; \pi)\|^{1+r} \leq \alpha, \quad \max_{0 \leq t \leq T-1} 2L_\pi \mu \|\Delta_t^{\pi_*}(\xi; \pi)\|^{1+r} + \|\Delta_t^{\pi_*}(\xi; \pi)\| \leq \eta, \quad (7)$$

then for all $1 \leq t \leq T$ the instantaneous imitation gap is bounded as

$$\|x_t^{\pi_*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} \mu \|\Delta_k^{\pi_*}(\xi; \pi)\|^{1+r}. \quad (8)$$

Theorem 3.1 shows that if a policy π is a sufficiently good approximation of the expert policy π_* on an expert trajectory $\{x_t^{\pi_*}(\xi)\}$, then the imitation gap $\|x_t^{\pi_*}(\xi) - x_t^\pi(\xi)\|$ can be upper bounded in terms of the discrepancy term $\max_{0 \leq k \leq t-1} \|\Delta_k^{\pi_*}(\xi; \pi)\|$ *evaluated on the expert trajectory*. To help illustrate the effect of the decay parameter r on the choices of μ and α , we make condition (6) more explicit by assuming that $\gamma(x) \leq Cx^{1+r}$ for all $x \in [0, 1]$. Then one can choose $\mu = 2^{1+r}C$ and $\alpha = \mathcal{O}(1)(L_\pi^{1+r}C)^{-1/r}$. This makes clear that a larger r , i.e., a more robust expert, leads to less restrictive conditions (7) on the policy π and a tighter upper bound on the imitation gap (8) (assuming $\max_{0 \leq k \leq t-1} \|\Delta_k^{\pi_*}(\xi; \pi)\| < 1$). In particular, for such systems, vanilla behavior cloning is sufficient to ensure bounded imitation gap. This also makes clear how the result breaks down when $r \approx 0$, i.e., when the decay rate is nearly linear, as the neighborhood α can become arbitrarily small, such that it may be impossible to learn a policy π that satisfies the bound (7) with a practical number of samples n . The interplay of the bounds (7) in Theorem 3.1 (and the subsequent theorems of its like) and the sample-complexity of imitation learning will be discussed in further detail in Section 4.

Slowly decaying class \mathcal{K} functions

When the class \mathcal{K} function $\gamma(x)$ decays to 0 slowly, for reasons discussed above, controlling the zeroth-order difference $\Delta_t^{\pi_*}(\xi; \pi)$ may not be sufficient to bound the imitation gap. In particular, we consider class \mathcal{K} functions satisfying $\gamma(x) \leq \mathcal{O}(x^{1/r})$ for some $r \geq 1$. Setting $p = \lfloor r \rfloor$, we now show that matching up to the p -th total derivative of π_* is sufficient to control the imitation gap. Analogously to Assumption 3.1, we make the following regularity assumption on the test policy π and expert policy π_* .

Assumption 3.2. *For a given non-negative $p \in \mathbb{N}$, assume that the test policy π and expert policy π_* are p -times continuously differentiable, and there exists a constant $L_{\partial^p \pi} \geq 0$ such that*

$$\|\bar{\pi}(x) - (J_{x_0}^p \bar{\pi})(x)\| \leq \frac{L_{\partial^p \pi}}{(p+1)!} \|x - x_0\|^{p+1}, \quad (9)$$

for all x, x_0 and $\bar{\pi} \in \{\pi, \pi_\}$, where $(J_{x_0}^p \bar{\pi})(x) := \sum_{j=0}^p \frac{1}{j!} (\partial_x^j \bar{\pi}(x_0))(x - x_0)^{\otimes j}$ is the p -th order Taylor polynomial of $\bar{\pi}$ evaluated at x_0 , and \otimes denotes the tensor product.*

With this assumption in hand, we provide the following guarantee on the imitation gap.

Theorem 3.2. *Let $f_{\text{cl}}^{\pi_\star}$ be η -locally δ -ISS for some $\eta > 0$, and assume that the class \mathcal{K} function $\gamma(\cdot)$ in (2) satisfies $\gamma(x) \leq \mathcal{O}(x^{1/r})$ for some $r \geq 1$. Fix a test policy π and initial condition $\xi \in \mathcal{X}$, and let Assumption 3.2 hold for $p \in \mathbb{N}$ satisfying $p + 1 - r > 0$. Choose $\mu, \alpha > 0$ such that*

$$2 \frac{L_{\partial^p \pi}}{(p+1)!} x^{p+1} + (x/\mu)^r \leq \gamma^{-1}(x), \text{ for all } 0 \leq x \leq \alpha \leq \frac{1}{2}. \quad (10)$$

Provided the j th total derivatives, $j = 0, \dots, p$, of the imitation error on the expert trajectory incurred by π satisfy:

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\| \right)^{1/r} \leq \alpha, \quad (11)$$

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\| \right)^{\frac{p+1}{r}} + \frac{2}{j!} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\| \leq \eta, \quad (12)$$

then for all $1 \leq t \leq T$ the instantaneous imitation gap is bounded by

$$\|x_t^{\pi_\star}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\| \right)^{1/r}. \quad (13)$$

Theorem 3.2 shows that if the p -th order Taylor series of the policy π approximately matches that of the expert policy π_\star when evaluated on an expert trajectory $\{x_t^{\pi_\star}(\xi)\}$, then the imitation gap $\|x_t^{\pi_\star}(\xi) - x_t^\pi(\xi)\|$ can be upper bounded in terms of the derivatives of the discrepancy term, i.e., by $\max_{0 \leq k \leq t-1} \max_{0 \leq j \leq p} \|\partial_x^j \Delta_k^{\pi_\star}(\xi; \pi)\|$, *evaluated on the expert trajectory*. To help illustrate the effect of the choice of the order p on the constants μ and α , we make condition (10) more explicit by assuming that $\gamma(x) \leq Cx^{1/r}$ for all $x \in [0, 1]$. Then one can choose $\mu = 2^{1/r}C$ and $\alpha = \mathcal{O}(1)(L_{\partial^p \pi} C^r)^{-1/(p+1-r)}$. This expression highlights a tradeoff: by picking larger order p , the right hand side α of bound (11) increases, but at the expense of having to match higher-order derivatives. This also highlights that both the order p and closeness required by Equation (11) get increasingly restrictive as r increases, matching our intuition that less robust experts lead to harder imitation learning problems.

Using estimated derivatives We show in Appendix C that the results of Theorems 3.1 and 3.2 extend gracefully to when only approximate derivatives $\widehat{\partial_x^j \pi_\star}(x)$ can be obtained, e.g., through finite-difference methods. In particular, if $\|\widehat{\partial_x^j \pi_\star}(x) - \partial_x^j \pi_\star(x)\| \leq \varepsilon$, then it suffices to appropriately tighten the constraints (11) and (12) by $\mathcal{O}(\varepsilon^{1/r})$ and $\mathcal{O}(\varepsilon)$, respectively. Please refer to Appendix C for more details.

4 Algorithms and generalization bounds for TaSIL

The analysis of Section 3 focused on a single test policy π and initial condition ξ . Theorems 3.1 and 3.2 motivate defining the p -TaSIL loss function:

$$\ell_p^{\pi_\star}(\xi; \pi) := \frac{1}{p+1} \sum_{j=0}^p \max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\|. \quad (14)$$

The corresponding policy $\hat{\pi}_{\text{TaSIL}, p}$ is the solution to the empirical risk minimization (ERM) problem:

$$\hat{\pi}_{\text{TaSIL}, p} \in \operatorname{argmin}_{\pi \in \Pi} \frac{1}{n} \sum_{i=1}^n \ell_p^{\pi_\star}(\xi_i; \pi), \quad (15)$$

which explicitly seeks to learn a policy $\pi \in \Pi$, for Π a suitable policy class, that matches the p -th order Taylor series expansion of the expert policy.⁵ In this section, we analyze the generalization and sample-complexity properties of the p -TaSIL ERM problem (15).

⁵Although we focus on the supremum loss $\max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\|$ in our analysis, we note that any surrogate loss that upper bounds the supremum loss, e.g., $\sum_{t=0}^{T-1} \|\partial_x^j \Delta_t^{\pi_\star}(\xi; \pi)\|$, can be used.

Our analysis in this section focuses on the *realizable* setting: we assume that for every dataset of expert trajectories $\{\{x_t^{\pi_*}(\xi_i)\}_{t=0}^{T-1}\}_{i=1}^n$, there exists a policy $\pi \in \Pi$ that achieves (near) zero empirical risk. This is true if, for example, $\pi_* \in \Pi$. In this setting, we demonstrate that we can attain generalization bounds that decay as $\tilde{\mathcal{O}}(n^{-1})$, where $\tilde{\mathcal{O}}(\cdot)$ hides poly-log dependencies on n . These rates are referred to as *fast* rates in statistical learning, since they decay faster than the $n^{-1/2}$ rate prescribed by the central limit theorem. We present analysis for the non-realizable setting in Appendix B: this analysis is standard and yields generalization bounds scaling as $\mathcal{O}(n^{-1/2})$.

Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a set of functions mapping some domain \mathcal{X} to $[0, 1]$.⁶ Let \mathcal{D} be a distribution with support restricted to \mathcal{X} , and denote the mean of $g \in \mathcal{G}$ with respect to $x \sim \mathcal{D}$ by $\mathbb{E}_x[g]$. Similarly, fixing data points $x_1, \dots, x_n \in \mathcal{X}$, we denote the empirical mean of g by $\mathbb{E}_n[g] := n^{-1} \sum_{i=1}^n g(x_i)$. We focus our analysis on the following class of parametric Lipschitz function classes.

Definition 4.1 (Lipschitz parametric function class). *A parametric function class $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ is called (B_θ, L_θ, q) -Lipschitz if $\mathcal{G} = \{g_\theta(\cdot) \mid \theta \in \Theta\}$ with $\Theta \subset \mathbb{R}^q$, and it satisfies the following boundedness and uniform Lipschitz conditions:*

$$\sup_{\theta \in \Theta} \|\theta\| \leq B_\theta, \quad \sup_{x \in \mathcal{X}} \sup_{\theta_1, \theta_2 \in \Theta, \theta_1 \neq \theta_2} \frac{|g_{\theta_1}(x) - g_{\theta_2}(x)|}{\|\theta_1 - \theta_2\|} \leq L_\theta. \quad (16)$$

We assume without loss of generality that $B_\theta L_\theta \geq 1$.

The description (16) is very general, and as we show next, is compatible with feed-forward neural networks with differentiable activation functions. We then have the following generalization bound, which adapts [23, Corollary 3.7] to Lipschitz parametric function classes using the machinery of local Rademacher complexities [23]. Alternatively, the result can also be derived from [24, Theorem 3].

Theorem 4.1. *Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a (B_θ, L_θ, q) -Lipschitz parametric function class. There exists a universal positive constant $K < 10^6$ such that the following holds. Given $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the i.i.d. draws $x_1, \dots, x_n \sim \mathcal{D}$, for all $g \in \mathcal{G}$, the following bound holds:*

$$\mathbb{E}_x[g] \leq 2\mathbb{E}_n[g] + K \left(\frac{q \log(B_\theta L_\theta n) + \log(1/\delta)}{n} \right). \quad (17)$$

We now use Theorem 4.1 to analyze the generalization properties of the p -TaSIL ERM problem (15). In what follows, we assume that the expert-closed loop system is stable in the sense of Lyapunov, i.e., that there exists $B_X > 0$ such that $\sup_{t \in \mathbb{N}, \xi \in \mathcal{X}} \|x_t^{\pi_*}(\xi)\| \leq B_X$, and consider the following parametric class of $p + 2$ continuously differentiable policies:

$$\Pi_{\theta,p} := \{\pi(x, \theta) \mid \theta \in \mathbb{R}^q, \|\theta\| \leq B_\theta, \pi(0, \theta) = 0 \forall \theta, \pi \text{ is } p + 2 \text{ continuously differentiable}\}. \quad (18)$$

Define the constants

$$B_j := \sup_{\|x\| \leq B_X, \|\theta\| \leq B_\theta} \|\partial_x^j \pi(x, \theta)\|, \quad L_j := \sup_{\|x\| \leq B_X, \|\theta\| \leq B_\theta} \|\partial_x^{j+1} \partial_\theta \pi(x, \theta)\|,$$

for $j = 0, \dots, p$, and note that they are guaranteed to be finite under our regularity assumptions. Finally, define the loss function class:

$$\ell_p^{\pi_*} \circ \Pi_{\theta,p} := \{\ell_p^{\pi_*}(\cdot; \pi) \text{ defined in (14)} \mid \pi \in \Pi_{\theta,p}\}. \quad (19)$$

From a repeated application of Taylor's theorem, we show in Lemma B.2 that $B_{\ell,p}^{-1}(\ell_p^{\pi_*} \circ \Pi_{\theta,p})$ is a $(B_\theta, B_{\ell,p}^{-1}L_{\ell,p}, q)$ -Lipschitz parametric function class for $B_{\ell,p} := \frac{2}{p+1} \sum_{j=0}^p B_j$ and $L_{\ell,p} := \frac{B_X}{p+1} \sum_{j=0}^p L_j$. We now combine this with Theorem 4.1 to bound the population risk achieved by the solution to the TaSIL ERM problem (15).

Corollary 4.1. *Let the policy class $\Pi_{\theta,p}$ be defined as in (18), and assume that $\pi_* \in \Pi_{\theta,p}$. Let the function class $\ell_p^{\pi_*} \circ \Pi_{\theta,p}$ be defined as in (19), and constants $B_{\ell,p}, L_{\ell,p}$ be defined as above. Let $\hat{\pi}_{\text{TaSIL},p}$ be any empirical risk minimizer (15). Then with probability at least $1 - \delta$ over the initial conditions $\{\xi_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$,*

$$\mathbb{E}_\xi[\ell_p^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},p})] \leq \mathcal{O}(1) B_{\ell,p} \frac{q \log(B_\theta B_{\ell,p}^{-1} L_{\ell,p} n) + \log(1/\delta)}{n}. \quad (20)$$

⁶This is without loss of generality for $[0, B]$ -bounded functions by considering the normalized function class $B^{-1}\mathcal{G} := \{B^{-1}g \mid g \in \mathcal{G}\} \subset [0, 1]^{\mathcal{X}}$.

We note that since these generalization bounds solely concern the supervised learning problem of matching the expert on the expert trajectories, the constants do not depend on the stability properties of the trajectories generated by the learned policy $\hat{\pi}_{\text{TaSIL},p}$. To convert the generalization bound (20) to a probabilistic bound on the imitation gap $\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},p})$, we first apply Markov's inequality to bound the probability that the conditions of Theorem 3.1 or 3.2 hold by the expected TaSIL loss (20), and then apply Corollary 4.1 together with Markov's inequality.

Theorem 4.2 (Rapidly decaying class \mathcal{K} functions). *Assume that $\pi_\star \in \Pi_{\theta,0}$ and let the assumptions of Theorem 3.1 hold for all $\pi \in \Pi_{\theta,0}$. Let Equation (6) hold with constants $\mu, \alpha > 0$, and assume without loss of generality that $\alpha/\mu \leq 1$, $L_\pi \mu \geq 1/2$. Let $\hat{\pi}_{\text{TaSIL},0}$ be an empirical risk minimizer of $\ell_0^{\pi_\star}$ over the policy class $\Pi_{\theta,0}$ for initial conditions $\{\xi_i\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$. Fix a failure probability $\delta \in (0, 1)$, and assume that*

$$n \geq \mathcal{O}(1) \max \left\{ B_{\ell,0} \frac{\kappa_\alpha}{\delta} \log \left(\frac{\kappa_\alpha B_{\theta} L_{\ell,0}}{\delta} \right), B_{\ell,0} \frac{\kappa_\eta}{\delta} \log \left(\frac{\kappa_\eta B_{\theta} L_{\ell,0}}{\delta} \right) \right\},$$

where $\kappa_\alpha := q(\mu/\alpha)^{1/(1+r)}$, $\kappa_\eta := qL_\pi \mu/\eta$. Then with probability at least $1 - \delta$, the imitation gap evaluated on $\xi \sim \mathcal{D}$ (drawn independently from $\{\xi_i\}_{i=1}^n$) satisfies

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},0}) \leq \mathcal{O}(1) \mu \left(\frac{1}{\delta} \frac{B_{\ell,0} q \log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n)}{n} \right)^{1+r}.$$

Theorem 4.3 (Slowly decaying class \mathcal{K} functions). *Assume that $\pi_\star \in \Pi_{\theta,p}$, and let the assumptions of Theorem 3.2 hold for all $\pi \in \Pi_{\theta,p}$. Let Equation (10) hold with constants $\mu, \alpha > 0$. Let $\hat{\pi}_{\text{TaSIL},p}$ be an empirical risk minimizer of $\ell_p^{\pi_\star}$ over the policy class $\Pi_{\theta,p}$ for initial conditions $\{\xi_i\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$. Fix a failure probability $\delta \in (0, 1)$, and assume*

$$n \geq \mathcal{O}(1) \max_{j \leq p} \max \left\{ B_j \frac{\kappa_{\alpha,j}}{\delta} \log \left(\frac{\kappa_{\alpha,j} B_{\theta} B_j^{-1} B_X L_j}{\delta} \right), B_j \frac{\kappa_{\eta,j}}{\delta} \log \left(\frac{\kappa_{\eta,j} B_{\theta} B_j^{-1} B_X L_j}{\delta} \right) \right\},$$

where $\kappa_{\alpha,j} := \left(\frac{\mu}{\alpha}\right)^r \frac{pq}{j!}$ and $\kappa_{\eta,j} := \left(\frac{L_{\partial p} \pi}{(p+1)!} \frac{\mu^{p+1}}{(j!)^{(p+1)/r}} + \frac{1}{j!}\right) \frac{pq}{\eta \delta}$. Then with probability at least $1 - \delta$, the imitation gap evaluated on $\xi \sim \mathcal{D}$ (drawn independently from $\{\xi_i\}_{i=1}^n$) satisfies

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},p}) \leq \mathcal{O}(1) \mu \max_{j \leq p} \left(\frac{p}{j! \delta} \frac{B_j q \log(B_{\theta} B_j^{-1} B_X L_j n)}{n} \right)^{1/r}.$$

In the rapidly decaying setting, corresponding to more robust experts, Theorem 4.2 states that $n \gtrsim \varepsilon^{-\frac{1}{1+r}}/\delta$ expert trajectories are sufficient to ensure that the imitation gap $\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},0}) \lesssim \varepsilon$ with probability at least $1 - \delta$. Recall that more robust experts have larger values of $r > 0$, leading to smaller sample-complexity bounds. In contrast, to achieve the same guarantees on the imitation gap in the slowly decaying setting, Theorem 4.3 states $n \gtrsim \varepsilon^{-r}/\delta$ expert trajectories are sufficient, where we recall that less robust experts have larger values of $r \geq 1$. These theorems quantitatively show how the robustness of an underlying expert affects the sample-complexity of IL, with more robust experts enjoying better dependence on ε than less robust experts. We note that analogous dependencies on expert stability are reflected in the burn-in requirements, i.e., the number of expert trajectories required to ensure with high probability no catastrophic distribution shift occurs, of each theorem.

5 Experiments

We compare three standard imitation learning algorithms, Behavior Cloning, DAgger, and DART, to TaSIL-augmented loss versions. In TaSIL-augmented algorithms, we replace the standard imitation loss function

$$\ell_{\text{IL}}^{\pi_\star}(\{\xi_i\}_{i=1}^n; \pi) := \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^{T-1} \|\Delta_t^{\pi_\star}(\xi_i; \pi)\|$$

with the p -TaSIL-augmented loss

$$\ell_{\text{TaSIL},p}^{\pi_\star}(\{\xi_i\}_{i=1}^n; \pi) := \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{T-1} \sum_{j=0}^p \lambda_j \|\text{vec}(\partial_x^j \Delta_t^{\pi_\star}(\xi_i; \pi))\|,$$

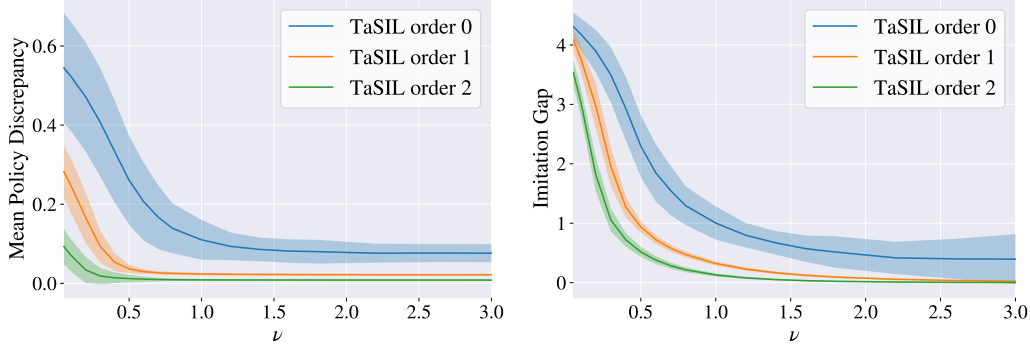


Figure 1: Left: The average Euclidean norm difference between the expert and learned policies on trajectories rolled out under the learned policy. **Right:** The maximum deviation between an expert and learned policy trajectory starting from identical initial conditions. All statistics are averaged across 50 test trajectories and we plot the mean and standard deviation for 10 random seeds.

where $\{\lambda_j\}_{j=0}^p$ are positive tunable regularization parameters. We use the Euclidean norm of the vectorized error in the derivative tensors as more optimizer-amenable surrogate to the operator norm.

We experimentally demonstrate (i) the effect of the expert stability properties and order of the TaSIL loss, and (ii) that TaSIL-loss-based imitation learning algorithms are more sample efficient than their standard counterparts. All experiments⁷ are carried out using Jax [25] GPU acceleration and automatic differentiation capabilities to compute the higher-order derivatives, and the Flax [26] neural network and Optax [27] optimization toolkits.

Stability Experiments To illustrate the effect of the expert closed-loop system stability on sample-complexity, we consider a simple δ -ISS stable dynamical system with a tunable γ input-to-state gain function. For state and input $x_t, u_t \in \mathbb{R}^{10}$, the dynamics are:

$$x_{t+1} = \eta x_t + (1 - \eta) \frac{\gamma(\|h(x_t) + u_t\|)}{\|h(x_t) + u_t\|} (h(x_t) + u_t).$$

The perturbation function $h : \mathbb{R}^{10} \rightarrow \mathbb{R}^{10}$ is set to a randomly initialized MLP with two hidden layers of width 32 and GELU [28] activations such that the expert $\pi_*(x) = -h(x)$ yield a closed loop system $f_{\text{cl}}^{\pi_*}(x, \Delta)$ which is δ -ISS stable with the specified class \mathcal{K} function γ (see Appendix D). We use $\eta = 0.95$ for all experiments presented here.

We investigate the performance of p -TaSIL loss functions for δ -ISS system with different class \mathcal{K} stability. We sweep \mathcal{K} functions $\gamma(x) = Cx^\nu$ for $\nu \in [0.05, 3]$, $C = 5$ and p -TaSIL loss functions for $p \in \{0, 1, 2\}$ (additional details can be found in Appendix D). The results of this sweep are shown in Figure 1. Higher-order p -TaSIL losses significantly reduce both the imitation gap and the mean policy discrepancy on test trajectories. Notably, the first and second order TaSIL loss maintain their improved performance for slower decaying class \mathcal{K} functions. Theorem 3.1 and Theorem 3.2 yield lower bounds of $\nu = 1$, $\nu = 2^{-1}$, and $\nu = 3^{-1}$ for closing the imitation gap using the 0-TaSIL, 1-TaSIL, and 2-TaSIL losses respectively. Figure 1 demonstrates significant performance degradation in policy discrepancy and decaying imitation gap starting around these threshold values.

MuJoCo Experiments We evaluate the ability of the TaSIL loss to improve performance on standard imitation learning tasks by modifying Behavior Cloning, DAgger [6], and DART [7] to use the $\ell_{\text{TaSIL},1}$ loss and testing them in simulation on different OpenAI Gym MuJoCo tasks [29]. The MuJoCo environments we use and their corresponding (state, input) dimensions are: Walker2d-v3 (17, 6), HalfCheetah-v3 (17, 6), Humanoid-v3 (376, 17), and Ant-v3 (111, 8).

For all environments we use pretrained expert policies obtained using Soft Actor Critic reinforcement learning by the Stable-Baselines3 [30] project. The experts consist of Multi-Layer Perceptrons with two hidden layers of 256 units each and ReLU activations. For all environments, learned policies have 2 hidden layers with 512 units each and GELU activations in addition to Batch Normalization. The final policy output for both the expert and learned policy are rescaled to the valid action space

⁷The code used for these experiments can be found at <https://github.com/unstable-zeros/TaSIL>

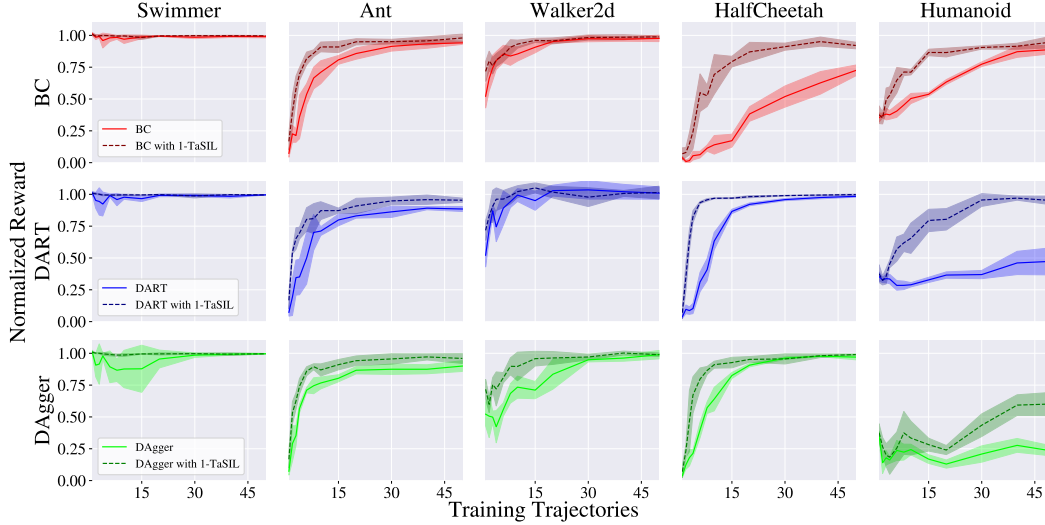


Figure 2: Cumulative expert-normalized rewards as a function of trajectory budget for policies trained using different algorithms with and without 1-TaSIL loss.

after applying a tanh nonlinearity. We used trajectories of length $T = 300$ for all experiments. We refer to Appendix E for additional experiment details.

In Figure 2 we report the mean expert-normalized rewards across 5 seeds for all algorithms and environments as a function of the trajectory budget provided. Algorithms with 1-TaSIL loss showed significant improvement in sample-complexity across all challenging environments. The expert for the Swimmer environment is very robust due to the simplicity of the task, and so as predicted by Theorem 3.1 and 4.2, all algorithms (with the exception of the vanilla DAgger algorithm due to it initially selecting poor rollout trajectories) are able to achieve near expert performance across all trajectory budgets. We are also able to nearly match or exceed the performance of standard on-policy methods DAgger and DART with our off-policy 1-TaSIL loss Behavior Cloning in all environments. Additional experimental results can be found in the appendix. These include representative videos of the behavior achieved by the expert, BC, and TaSIL-augmented BC policies in the supplementary material, where once again, a striking improvement is observed, especially in low-data regimes for harder environments, as well as a systematic study of finite-difference-based approximations of 1-TaSIL which achieve comparable performance to Jacobian-based implementations.

6 Conclusion

We presented Taylor Series Imitation Learning (TaSIL), a simple augmentation to behavior cloning that penalizes deviations in the higher-order Taylor series terms between the learned and expert policies. We showed that δ -ISS experts are easier to learn, both in terms of the loss-function that needs to be optimized and sample-complexity guarantees. Finally, we showed the benefit of using TaSIL-augmented losses in BC, DAgger, and DART across a variety of MuJoCo tasks. This work opens up many exciting future directions, including extending TaSIL to pixel-based IL, and to (offline/inverse) reinforcement learning settings.

Acknowledgements

We thank Vikas Sindhwani, Sumeet Singh, Jean-Jacques E. Slotine, Jake Varley, and Fengjun Yang for helpful feedback. Nikolai Matni is supported by NSF awards CPS-2038873, CAREER award ECCS-2045834, and a Google Research Scholar award.

References

- [1] Ahmed Hussein, Mohamed M. Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- [2] Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J. Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning. *Foundations and Trends® in Robotics*, 7(1–2):1–179, 2018.
- [3] Dean A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems*, volume 1, 1988.
- [4] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4693–4700, 2018.
- [5] Stefan Schaal. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3(6):233–242, 1999.
- [6] Stéphane Ross, Geoffrey Gordon, and Drew 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*, volume 15, pages 627–635. PMLR, 2011.
- [7] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning. In *Proceedings of the 1st Annual Conference on Robot Learning*, volume 78, pages 143–156. PMLR, 2017.
- [8] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [9] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [10] Michael Hertneck, Johannes Köhler, Sebastian Trimpe, and Frank Allgöwer. Learning an approximate model predictive controller with guarantees. *IEEE Control Systems Letters*, 2(3): 543–548, 2018.
- [11] He Yin, Peter Seiler, Ming Jin, and Murat Arcak. Imitation learning with stability and safety guarantees. *IEEE Control Systems Letters*, 6:409–414, 2022.
- [12] Allen Ren, Sushant Veer, and Anirudha Majumdar. Generalization guarantees for imitation learning. In *Proceedings of the 2020 Conference on Robot Learning*, volume 155, pages 1426–1442. PMLR, 2021.
- [13] Stephen Tu, Alexander Robey, Tingnan Zhang, and Nikolai Matni. On the sample complexity of stability constrained imitation learning. *arXiv preprint arXiv:2102.09161*, 2021.
- [14] Winfried Lohmiller and Jean-Jacques E. Slotine. On contraction analysis for non-linear systems. *Automatica*, 34(6):683–696, 1998.
- [15] Sumeet Singh, Spencer M. Richards, Vikas Sindhvani, Jean-Jacques E. Slotine, and Marco Pavone. Learning stabilizable nonlinear dynamics with contraction-based regularization. *The International Journal of Robotics Research*, 40:1123–1150, 2020.
- [16] Andre Lemme, Klaus Neumann, R. Felix Reinhart, and Jochen J. Steil. Neural learning of vector fields for encoding stable dynamical systems. *Neurocomputing*, 141:3–14, 2014.
- [17] Harish Ravichandar, Iman Salehi, and Ashwin Dani. Learning partially contracting dynamical systems from demonstrations. In *Proceedings of the 1st Annual Conference on Robot Learning*, volume 78, pages 369–378. PMLR, 2017.
- [18] Vikas Sindhvani, Stephen Tu, and Mohi Khansari. Learning contracting vector fields for stable imitation learning. *arXiv preprint arXiv:1804.04878*, 2018.

- [19] Nicholas M. Boffi, Stephen Tu, and Jean-Jacques E. Slotine. Regret bounds for adaptive nonlinear control. In *Proceedings of the 3rd Conference on Learning for Dynamics and Control*, volume 144, pages 471–483. PMLR, 2021.
- [20] Nicholas M. Boffi, Stephen Tu, Nikolai Matni, Jean-Jacques E. Slotine, and Vikas Sindhwani. Learning stability certificates from data. In *Proceedings of the 2020 Conference on Robot Learning*, volume 155, pages 1341–1350. PMLR, 2021.
- [21] David Angeli. A lyapunov approach to incremental stability properties. *IEEE Transactions on Automatic Control*, 47(3):410–421, 2002.
- [22] Hassan K. Khalil. *Nonlinear Systems*. Pearson Education. Prentice Hall, 2002.
- [23] Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- [24] David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. *Information and Computation*, 100(1):78–150, 1992.
- [25] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- [26] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. URL <http://github.com/google/flax>.
- [27] Matteo Hessel, David Budden, Fabio Viola, Mihaela Rosca, Eren Sezener, and Tom Hennigan. Optax: composable gradient transformation and optimisation, in jax!, 2020. URL <http://github.com/deepmind/optax>.
- [28] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [29] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [30] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021.
- [31] Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, volume 23, 2010.
- [32] Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- [33] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [34] Olivier Bousquet. *Concentration inequalities and empirical processes theory applied to the analysis of learning algorithms*. PhD thesis, École Polytechnique: Department of Applied Mathematics Paris, France, 2002.
- [35] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- [36] Max Simchowitz, Horia Mania, Stephen Tu, Michael I. Jordan, and Benjamin Recht. Learning without mixing: Towards a sharp analysis of linear system identification. In *Proceedings of the 31st Conference On Learning Theory*, volume 75, pages 439–473. PMLR, 2018.

- [37] Franco Woolfe, Edo Liberty, Vladimir Rokhlin, and Mark Tygert. A fast randomized algorithm for the approximation of matrices. *Applied and Computational Harmonic Analysis*, 25(3): 335–366, 2008.
- [38] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 2012.

Contents

1	Introduction	1
1.1	Related work	2
2	Problem formulation	3
3	Bounding the imitation gap on a single trajectory	3
4	Algorithms and generalization bounds for TaSIL	6
5	Experiments	8
6	Conclusion	10
A	Proofs for Section 3	15
B	Proofs for Section 4	18
B.1	Preliminaries	18
B.2	Generalization bound for the non-realizable setting	18
B.3	Proof of Theorem 4.1	20
B.4	Proof of Corollary 4.1	22
B.5	Proofs of Theorem 4.2 and Theorem 4.3	23
C	Using finite-differencing to approximate derivatives	26
C.1	Satisfying Conditions (11) and (12) with approximate derivatives	26
C.2	Practical approaches for approximating derivatives	26
C.3	Experimental results	27
D	Additional information for stability experiments	28
E	Additional information for MuJoCo experiments	28

A Proofs for Section 3

Proposition 3.1. *Let the expert closed-loop system $f_{\text{cl}}^{\pi^*}$ be η -locally δ -ISS for some $\eta > 0$. Fix an imitation gap bound $\varepsilon > 0$, initial condition ξ , and policy π . Then if*

$$\max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_*(x_t^{\pi^*}(\xi) + \delta) - \pi(x_t^{\pi^*}(\xi) + \delta)\| \leq \min\{\eta, \gamma^{-1}(\varepsilon)\}, \quad (5)$$

we have that the imitation gap satisfies $\Gamma_T(\xi; \pi) \leq \varepsilon$.

Proof. We do a proof by induction.

Base case $t = 0$: We trivially have at $t = 0$:

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| = \|\xi - \xi\| = 0 \leq \varepsilon.$$

Induction step: Assume for some $k > 0$, we have $\max_{t \leq k-1} \|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \varepsilon$. We set $\delta := x_{k-1}^{\pi^*}(\xi) - x_{k-1}^\pi(\xi)$ such that $\|\delta\| \leq \varepsilon$. From Equation (5), we are guaranteed that

$$\begin{aligned} \|\pi_*(x_{k-1}^{\pi^*}(\xi)) - \pi(x_{k-1}^{\pi^*}(\xi))\| &= \|\pi_*(x_{k-1}^{\pi^*}(\xi) + \delta) - \pi(x_{k-1}^{\pi^*}(\xi) + \delta)\| \\ &\leq \max_{0 \leq t \leq k-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_*(x_t^{\pi^*}(\xi) + \delta) - \pi(x_t^{\pi^*}(\xi) + \delta)\| \\ &\leq \min\{\eta, \gamma^{-1}(\varepsilon)\}. \end{aligned}$$

Since $f_{\text{cl}}^{\pi^*}$ is η -locally δ -ISS, we get from (3)

$$\begin{aligned} \|x_k^\pi(\xi) - x_k^{\pi^*}(\xi)\| &\leq \gamma \left(\max_{0 \leq s \leq k-1} \|\pi_*(x_s^{\pi^*}(\xi)) - \pi(x_s^{\pi^*}(\xi))\| \right) \\ &\leq \gamma(\min\{\eta, \gamma^{-1}(\varepsilon)\}) \\ &\leq \varepsilon, \end{aligned}$$

and thus $\max_{t \leq k} \|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \varepsilon$, completing the induction step. \square

Theorem 3.1. *Fix a test policy π and initial condition $\xi \in \mathcal{X}$, and let Assumption 3.1 hold. Let $f_{\text{cl}}^{\pi^*}$ be η -locally δ -ISS for some $\eta > 0$, and assume that the class \mathcal{K} function $\gamma(\cdot)$ in (2) satisfies $\gamma(x) \leq \mathcal{O}(x^{1+r})$ for some $r > 0$. Choose constants $\mu, \alpha > 0$ such that*

$$2L_\pi x + (x/\mu)^{\frac{1}{1+r}} \leq \gamma^{-1}(x) \text{ for all } 0 \leq x \leq \alpha. \quad (6)$$

Provided that the imitation error on the expert trajectory incurred by π satisfies:

$$\max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r} \leq \alpha, \quad \max_{0 \leq t \leq T-1} 2L_\pi \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r} + \|\Delta_t^{\pi^*}(\xi; \pi)\| \leq \eta, \quad (7)$$

then for all $1 \leq t \leq T$ the instantaneous imitation gap is bounded as

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} \mu \|\Delta_k^{\pi^*}(\xi; \pi)\|^{1+r}. \quad (8)$$

Proof. In order to leverage Proposition 3.1 we must first find a solution ε to Equation (5). By Lipschitzness of the policy class,

$$\max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_*(x_t^{\pi^*}(\xi) + \delta) - \pi(x_t^{\pi^*}(\xi) + \delta)\| \leq 2L_\pi \varepsilon + \max_{0 \leq t \leq T-1} \|\Delta_t^{\pi^*}(\xi; \pi)\|,$$

and using the lower bound in Equation (6) it is therefore sufficient to find a solution $\varepsilon \leq \alpha$ to

$$\begin{aligned} 2L_\pi \varepsilon + \max_{0 \leq t \leq T-1} \|\Delta_t^{\pi^*}(\xi; \pi)\| &\leq 2L_\pi \varepsilon + (\varepsilon/\mu)^{\frac{1}{1+r}} \\ \iff \max_{0 \leq t \leq T-1} \|\Delta_t^{\pi^*}(\xi; \pi)\| &\leq (\varepsilon/\mu)^{\frac{1}{1+r}}. \end{aligned}$$

Picking $\varepsilon = \max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r}$ and adding the constraint $\varepsilon \leq \alpha$ in order to ensure the solution is sufficiently small allows use to apply Proposition 3.1 and obtain the final result

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r}.$$

Provided that

$$\max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r} \leq \alpha, \quad \max_{0 \leq t \leq T-1} 2L_\pi \mu \|\Delta_t^{\pi^*}(\xi; \pi)\|^{1+r} + \|\Delta_t^{\pi^*}(\xi; \pi)\| \leq \eta$$

Thus completing the proof. \square

Theorem 3.2. Let $f_{\text{cl}}^{\pi^*}$ be η -locally δ -ISS for some $\eta > 0$, and assume that the class \mathcal{K} function $\gamma(\cdot)$ in (2) satisfies $\gamma(x) \leq \mathcal{O}(x^{1/r})$ for some $r \geq 1$. Fix a test policy π and initial condition $\xi \in \mathcal{X}$, and let Assumption 3.2 hold for $p \in \mathbb{N}$ satisfying $p + 1 - r > 0$. Choose $\mu, \alpha > 0$ such that

$$2 \frac{L_{\partial^p \pi}}{(p+1)!} x^{p+1} + (x/\mu)^r \leq \gamma^{-1}(x), \text{ for all } 0 \leq x \leq \alpha \leq \frac{1}{2}. \quad (10)$$

Provided the j th total derivatives, $j = 0, \dots, p$, of the imitation error on the expert trajectory incurred by π satisfy:

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{1/r} \leq \alpha, \quad (11)$$

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{\frac{p+1}{r}} + \frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \leq \eta, \quad (12)$$

then for all $1 \leq t \leq T$ the instantaneous imitation gap is bounded by

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{1/r}. \quad (13)$$

Proof. We proceed similarly as in the proof of Theorem 3.1. From Proposition 3.1, we can take the p th Taylor expansion of the left hand side of Equation (5) and apply the triangle inequality a few times to yield:

$$\begin{aligned} & \max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_\star(x_t^{\pi^*}(\xi) + \delta) - \pi(x_t^{\pi^*}(\xi) + \delta)\| \\ & \leq \max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_\star(x_t^{\pi^*}(\xi)) - \pi(x_t^{\pi^*}(\xi))\| \\ & \quad + \|\pi_\star(x_t^{\pi^*}(\xi) + \delta) - \pi_\star(x_t^{\pi^*}(\xi)) - (\pi(x_t^{\pi^*}(\xi) + \delta) - \pi(x_t^{\pi^*}(\xi)))\| \\ & \leq \max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \|\pi_\star(x_t^{\pi^*}(\xi)) - \pi(x_t^{\pi^*}(\xi))\| \\ & \quad + \left\| \sum_{j=1}^p \frac{1}{j!} \partial_x^j \pi_\star(x_t^{\pi^*}(\xi)) \cdot \delta^{\otimes j} - \sum_{j=1}^p \frac{1}{j!} \partial_x^j \pi(x_t^{\pi^*}(\xi)) \cdot \delta^{\otimes j} \right\| + 2 \frac{L_{\partial^p \pi}}{(p+1)!} \|\delta\|^{p+1} \\ & \leq \max_{0 \leq t \leq T-1} \sup_{\|\delta\| \leq \varepsilon} \sum_{j=0}^p \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \cdot \delta^{\otimes j}\| + 2 \frac{L_{\partial^p \pi}}{(p+1)!} \|\delta\|^{p+1} \\ & \leq \max_{0 \leq t \leq T-1} \sum_{j=0}^p \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \varepsilon^j + 2 \frac{L_{\partial^p \pi}}{(p+1)!} \varepsilon^{p+1}. \end{aligned}$$

Therefore, it suffices to find an ε small enough such that

$$\max_{0 \leq t \leq T-1} 2 \frac{L_{\partial^p \pi}}{(p+1)!} \varepsilon^{p+1} + \sum_{j=0}^p \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \varepsilon^j \leq \gamma^{-1}(\varepsilon).$$

Since we are given $\gamma(x) \leq \mathcal{O}(x^{1/r})$, we have $\gamma^{-1}(x) \geq \Omega(x^r)$. This motivates finding a large enough μ and small enough neighborhood α such that

$$\max_{0 \leq t \leq T-1} 2 \frac{L_{\partial^p \pi}}{(p+1)!} \varepsilon^{p+1} + \left(\frac{\varepsilon}{\mu} \right)^r \leq \gamma^{-1}(\varepsilon),$$

for all $0 < \varepsilon \leq \alpha \leq 1/2$. In essence, we want to find a sufficiently small neighborhood α such that the ε^{p+1} term is dominated by the ε^r term, while also selecting a μ such that the total sum is still upper bounded by $\gamma^{-1}(x) \geq \Omega(x^r)$ in this neighborhood. The choice of raising ε/μ to the r -th power arises from the fact that r is the smallest exponent—thus affecting the imitation gap in Equation (13) downstream least severely—that ensures μ, α will always exist. Having found such μ, α ,

we now simply have to find $\|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\|$ small enough such that

$$\begin{aligned} \sum_{j=0}^p \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \varepsilon^j &\leq \max_{j \leq p} \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \sum_{j=0}^p \varepsilon^j \\ &\leq \max_{j \leq p} \frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \quad \varepsilon \leq \alpha \leq 1/2 \\ &= \left(\frac{\varepsilon}{\mu}\right)^r. \end{aligned} \quad (21)$$

Solving this for ε , we get

$$\varepsilon = \max_{j \leq p} \mu \left(\frac{2}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{1/r},$$

as long as $\varepsilon \leq \alpha$, the neighborhood condition, and $2 \frac{L_{\partial^p \pi}}{(p+1)!} \varepsilon^{p+1} + \left(\frac{\varepsilon}{\mu}\right)^r \leq \eta$, the locality for δ -ISS. These correspond to the conditions (11) and (12), respectively. This completes the proof. \square

For completeness we present here a stronger variant of Theorem 3.2 for the special case where $p = r \in \mathbb{N}$. In this scenario we are able to remove the dependency of the imitation gap bounds on the p th order derivative provided it can be made sufficiently small.

Theorem A.1. *Let $f_{\text{cl}}^{\pi^*}$ be η -locally δ -ISS for some $\eta > 0$, and assume that the class \mathcal{K} function $\gamma(\cdot)$ in (2) satisfies $\gamma(x) \leq \mathcal{O}(x^{1/r})$ for some $r \geq 1$. Fix a test policy π and initial condition $\xi \in \mathcal{X}$, and let Assumption 3.2 hold with $p = r \in \mathbb{N}$. Choose $\mu, \alpha > 0$ such that*

$$2 \frac{L_{\partial^p \pi}}{(p+1)!} x^{p+1} + (x/\mu)^p \leq \gamma^{-1}(x), \text{ for all } 0 \leq x \leq \alpha \leq \frac{1}{2}. \quad (22)$$

Provided the j th total derivatives, $j = 0, \dots, p$, of the imitation error on the expert trajectory incurred by π satisfy:

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p-1} \mu \left(\frac{4}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{1/p} \leq \alpha, \quad (23)$$

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p-1} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{4}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{\frac{p+1}{p}} + \frac{4}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \leq \eta, \quad (24)$$

$$\|\partial_x^p \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{p!}{2\mu^p} \quad (25)$$

then for all $1 \leq t \leq T$ the instantaneous imitation gap is bounded by

$$\|x_t^{\pi^*}(\xi) - x_t^{\pi}(\xi)\| \leq \max_{0 \leq k \leq t-1} \max_{0 \leq j \leq p-1} \mu \left(\frac{4}{j!} \right)^{1/r} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\|^{1/r}. \quad (26)$$

Proof. We follow the proof of Theorem 3.2 until Equation (21). We then wish to solve

$$\sum_{j=0}^p \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \varepsilon^j \leq \left(\frac{\varepsilon}{\mu}\right)^p.$$

Since the order of the RHS is p , provided that $\frac{1}{p!} \|\partial_x^p \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{1}{2} \frac{1}{\mu^p}$ we can write

$$\sum_{j=0}^{p-1} \frac{1}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \varepsilon^j \leq \frac{1}{2} \left(\frac{\varepsilon}{\mu}\right)^p.$$

Upper-bounding the polynomial on the LHS using a geometric series and solving for ε we get

$$\varepsilon = \max_{j \leq p-1} \mu \left(\frac{4}{j!} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \right)^{1/p},$$

provided that $\varepsilon \leq \alpha$, $2 \frac{L_{\partial^p \pi}}{(p+1)!} \varepsilon^{p+1} + \left(\frac{\varepsilon}{\mu}\right)^p \leq \eta$, and $\|\partial_x^p \Delta_t^{\pi^*}(\xi; \pi)\| < \frac{p!}{2\mu^p}$. These conditions correspond to that of the theorem, completing the proof. \square

Corollary A.1. Consider a δ -ISS $f_{\text{cl}}^{\pi^*}$ system with $\gamma(x) := \gamma x, \gamma > 0$ and $\eta = \infty$. Let Assumption 3.2 hold with $p = 1$ and assume without loss of generality $\gamma L_{\partial\pi} \geq 1$. Provided

$$\max_{0 \leq t \leq T-1} \|\partial_x \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{1}{4\gamma}, \quad \max_{0 \leq t \leq T-1} \|\Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{1}{16\gamma^2 L_{\partial\pi}}$$

then for all $0 \leq t \leq T$

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} 8\gamma \|\Delta_k^{\pi^*}(\xi; \pi)\|.$$

Proof. Choose $\alpha := \frac{1}{2\gamma L_{\partial\pi}}$ and $\mu := 2\gamma$. Assume $\gamma L_{\partial\pi} \geq 1$. Since $\gamma^{-1}(x) = \frac{x}{\gamma}$, for $x \leq \alpha$ it holds that

$$L_{\partial^p \pi} x^2 + (x/\mu) \leq \frac{x}{2\gamma} + \frac{x}{2\gamma} \leq \gamma^{-1}(x) := \frac{x}{\gamma}.$$

and we can directly apply the $p = r = 1$ special case of Theorem A.1. Then, if the constraints described by Equations (23) and (25) are satisfied:

$$\max_{0 \leq t \leq T-1} \|\partial_x \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{p!}{2\mu^p} = \frac{1}{4\gamma}, \quad \max_{0 \leq t \leq T-1} \|\Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{0!}{4} \left(\frac{\alpha}{\mu}\right)^p = \frac{1}{16\gamma^2 L_{\partial\pi}},$$

it holds for all $1 \leq t \leq T$

$$\|x_t^{\pi^*}(\xi) - x_t^\pi(\xi)\| \leq \max_{0 \leq k \leq t-1} 8\gamma \|\Delta_k^{\pi^*}(\xi; \pi)\|.$$

□

B Proofs for Section 4

B.1 Preliminaries

Let $\mathcal{G} \subset \mathbb{R}^{\mathcal{X}}$ be a set of functions, and let $x_1, \dots, x_n \in \mathcal{X}$ be a fixed set of points. We will endow \mathcal{G} with the following empirical L^2 pseudo-metric space structure:

$$d(f, g) := \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - g(x_i))^2}, \quad f, g \in \mathcal{G}.$$

The empirical Rademacher complexity of \mathcal{G} is defined as:

$$\mathcal{R}_n(\mathcal{G}) := \mathbb{E}_\varepsilon \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i g(x_i) \right],$$

where the $\{\varepsilon_i\}_{i=1}^n$ are independent Rademacher random variables. Dudley's inequality yields a bound on $\mathcal{R}_n(\mathcal{G})$ using the metric space structure of (\mathcal{G}, d) .

Lemma B.1 (Dudley's inequality [cf. 31, Lemma A.3]). *Let $R := \sup_{f \in \mathcal{G}} d(f, 0)$ be the radius of the set \mathcal{G} . We have that:*

$$\mathcal{R}_n(\mathcal{G}) \leq \inf_{\alpha \in [0, R]} \left\{ 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^R \sqrt{\log N(\mathcal{G}; d, \varepsilon)} d\varepsilon \right\}.$$

Here, $N(\mathcal{G}; d, \varepsilon)$ denotes the covering number of \mathcal{G} in the metric d at resolution ε .

B.2 Generalization bound for the non-realizable setting

We use standard techniques to derive a generalization bound for the *non-realizable setting*, i.e., where π_* may not necessarily be contained in the hypothesis class Π . Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a given function class. We have the following standard uniform convergence generalization bound [cf. 32, Theorem 4.10]: with probability greater than $1 - \delta$ over $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$, we have

$$\sup_{g \in \mathcal{G}} |\mathbb{E}_x[g] - \mathbb{E}_n[g]| \leq 2\mathbb{E}_{x_{1:n}}[\mathcal{R}_n(\mathcal{G})] + \sqrt{\frac{\log(2/\delta)}{n}}, \quad (27)$$

where $\mathbb{E}_{x_{1:n}}$ denotes expectation over the randomness of x_1, \dots, x_n . To establish an upper bound on $\mathbb{E}_{x_{1:n}}[\mathcal{R}_n(\mathcal{G})]$, we focus on the Lipschitz parametric case, though we note many analogous bounds can be computed for a plethora of other function classes [32].

Theorem B.1. Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a (B_θ, L_θ, q) -Lipschitz parametric function class. Given $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the i.i.d. draws $x_1, \dots, x_n \sim \mathcal{D}$, the following bound holds:

$$\sup_{g \in \mathcal{G}} |\mathbb{E}_x[g] - \mathbb{E}_n[g]| \leq 48 \sqrt{\frac{q \log(3B_\theta L_\theta)}{n}} + \sqrt{\frac{\log(2/\delta)}{n}}. \quad (28)$$

Proof. This argument is fairly standard. Fix a set of points $x_1, \dots, x_n \in \mathcal{X}$. Since \mathcal{G} contains only functions with range $[0, 1]$, the radius of the set \mathcal{G} in the empirical L^2 metric is:

$$\sup_{f \in \mathcal{G}} d(f, 0) \leq 1.$$

Therefore, Dudley's inequality (Lemma B.1) yields:

$$\mathcal{R}_n(\mathcal{G}) \leq \frac{12}{\sqrt{n}} \int_0^1 \sqrt{\log N(\mathcal{G}; d, \varepsilon)} d\varepsilon.$$

Now using the fact that \mathcal{G} is a (B_θ, L_θ, q) -Lipschitz parametric function class, it is not hard to see that for any $\varepsilon > 0$, an $\varepsilon/(B_\theta L_\theta)$ -cover of $\mathbb{B}_2^q(1)$ in the Euclidean metric yields an ε -cover of \mathcal{G} in the d -metric. Hence, for any $\varepsilon \in (0, 1)$, by a standard volume comparison argument:

$$\begin{aligned} \log N(\mathcal{G}; d, \varepsilon) &\leq \log N\left(\mathbb{B}_2^q(1); \|\cdot\|, \frac{\varepsilon}{B_\theta L_\theta}\right) \\ &\leq q \log\left(1 + \frac{2B_\theta L_\theta}{\varepsilon}\right) \\ &\leq q \log\left(\frac{3B_\theta L_\theta}{\varepsilon}\right). \end{aligned}$$

Therefore, we have:

$$\begin{aligned} \int_0^1 \sqrt{\log N(\mathcal{G}; d, \varepsilon)} d\varepsilon &\leq \sqrt{q} \int_0^1 \sqrt{\log\left(\frac{3B_\theta L_\theta}{\varepsilon}\right)} d\varepsilon \\ &\leq \sqrt{q \log(3B_\theta L_\theta)} + \sqrt{q} \int_0^1 \sqrt{\log(1/\varepsilon)} d\varepsilon \quad \text{using } \sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \\ &\leq \sqrt{q \log(3B_\theta L_\theta)} + \sqrt{q} \quad \text{using } \int_0^1 \sqrt{\log\left(\frac{1}{\varepsilon}\right)} d\varepsilon \leq 1 \\ &\leq 2\sqrt{q \log(3B_\theta L_\theta)}. \end{aligned}$$

Plugging this back into Dudley's inequality:

$$\mathcal{R}_n(\mathcal{G}) \leq 24 \sqrt{\frac{q}{n}} \sqrt{\log(3B_\theta L_\theta)}.$$

The claim now follows from the standard uniform convergence inequality (27). \square

Applying this generalization bound to the $(B_\theta, B_{\ell,p}^{-1} L_{\ell,p}, q)$ -Lipschitz parametric function class $B_{\ell,p}^{-1}(\ell_p^{\pi^*} \circ \Pi_{\theta,p})$, we get the non-realizable analogue to Corollary 4.1.

Corollary B.1. Let the policy class $\Pi_{\theta,p}$ be defined as in (18). Let the function class $\ell_p^{\pi^*} \circ \Pi_{\theta,p}$ be defined as in (19), and constants $B_{\ell,p}, L_{\ell,p}$ be defined as above. Let $\hat{\pi}_{\text{TaSIL},p}$ be any empirical risk minimizer (15). Then with probability at least $1 - \delta$ over the initial conditions $\{\xi_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$,

$$\mathbb{E}_\xi[\ell_p^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},p})] \leq \mathbb{E}_n[\ell_p^{\pi^*}(\cdot; \hat{\pi}_{\text{TaSIL},p})] + 48 B_{\ell,p} \sqrt{\frac{q \log(3B_\theta B_{\ell,p}^{-1} L_{\ell,p})}{n}} + B_{\ell,p} \sqrt{\frac{\log(2/\delta)}{n}}. \quad (29)$$

Inserting the generalization bound in Corollary B.1 in lieu of Corollary 4.1 for the rest of the bounds seen in Section 4 yields the sample complexity bounds relevant to our problem in the non-realizable setting. However, we note an important subtlety that manifests in the non-realizable regime. We note that in Corollary 4.1, due to realizability, the generalization bound monotonically decreases to 0 with n , whereas in Corollary B.1, we have an additive factor of $\mathbb{E}_n[\ell_p^{\pi^*}(\cdot; \hat{\pi}_{\text{TASIL},p})]$. It is therefore possible for either small enough n or insufficiently expressive function classes $\Pi_{\theta,p}$ that the non-zero empirical risk automatically violates the imitation error requirements in Theorems 3.1 and 3.2. Thus, a necessary assumption must be made in the non-realizable setting for the function class to be expressive enough such that the empirical risk it incurs on sufficiently large datasets satisfies the imitation error requirements with high probability.

B.3 Proof of Theorem 4.1

Before turning to the proof of Theorem 4.1, we introduce some notation and tools from the local Rademacher complexity literature [33, 34].

Definition B.1 (Sub-root function). *A function $\phi : [0, \infty) \rightarrow \mathbb{R}$ is said to be a sub-root function if:*

- a) ϕ is non-negative.
- b) ϕ is not the zero function.
- c) ϕ is non-decreasing.
- d) $r \mapsto \phi(r)/\sqrt{r}$ is non-increasing.

For any non-negative function class \mathcal{G} , scalar $r \geq 0$, and n points $x_1, \dots, x_n \in \mathcal{X}$, define:

$$\mathcal{H}_n(r; x_{1:n}) := \{g \in \mathcal{G} \mid \mathbb{E}_n[g] \leq r\}.$$

The following is from Bousquet [34].

Theorem B.2 (Bousquet [34, Theorem 6.1]). *Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$, and fix a $\delta \in (0, 1)$. With probability at least $1 - \delta$ over the i.i.d. draws of x_1, \dots, x_n , the following holds. Let ϕ_n be any sub-root function (cf. Definition B.1) satisfying:*

$$\mathcal{R}_n(\mathcal{H}_n(r; x_{1:n})) \leq \phi_n(r), \quad \forall r > 0.$$

Let r_n^ denote the largest solution to the equation $\phi_n(r) = r$. Then, for all $g \in \mathcal{G}$:*

$$\mathbb{E}_x[g] \leq 2\mathbb{E}_n[g] + 106r_n^* + \frac{48(\log(1/\delta) + 6 \log \log n)}{n}.$$

With these definitions and preliminary results in place, we turn to the proof of Theorem 4.1.

Theorem 4.1. *Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a (B_θ, L_θ, q) -Lipschitz parametric function class. There exists a universal positive constant $K < 10^6$ such that the following holds. Given $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the i.i.d. draws $x_1, \dots, x_n \sim \mathcal{D}$, for all $g \in \mathcal{G}$, the following bound holds:*

$$\mathbb{E}_x[g] \leq 2\mathbb{E}_n[g] + K \left(\frac{q \log(B_\theta L_\theta n) + \log(1/\delta)}{n} \right). \quad (17)$$

Proof. Fix a set of points $x_1, \dots, x_n \in \mathcal{X}$. Define $\mathcal{G}_n(r; x_{1:n})$ as:

$$\mathcal{G}_n(r; x_{1:n}) := \{g \in \mathcal{G} \mid \mathbb{E}_n[g^2] \leq r\}.$$

For what follows, we often suppress the explicit dependence on $x_{1:n}$ in the notation for \mathcal{H}_n and \mathcal{G}_n . Observe that since $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$, we have $\mathbb{E}_n[g^2] \leq \mathbb{E}_n[g]$ for every $g \in \mathcal{G}$, and therefore:

$$\mathcal{H}_n(r) \subseteq \mathcal{G}_n(r), \quad \forall r \geq 0.$$

Hence $\mathcal{R}_n(\mathcal{H}_n(r)) \leq \mathcal{R}_n(\mathcal{G}_n(r))$, and it suffices for us to prove an upper bound on the latter.

Proposition B.1. *Let $\mathcal{G} \subset [0, 1]^{\mathcal{X}}$ be a (B_θ, L_θ, q) -Lipschitz parametric function class. Fix a set of points $x_1, \dots, x_n \in \mathcal{X}$. We have that:*

$$\mathcal{R}_n(\mathcal{G}_n(r; x_{1:n})) \leq 24\sqrt{2} \sqrt{\frac{q}{n}} \min\{\sqrt{r}, 1\} \sqrt{\log \left(\frac{6B_\theta L_\theta}{\min\{\sqrt{r}, 1\}} \right)}.$$

Proof of Proposition B.1. The radius of the set $\mathcal{G}_n(r)$ in the empirical L^2 metric d is upper bounded by \sqrt{r} by definition. Furthermore, the radius of \mathcal{G} in the metric d is upper bounded by one. Hence, since $\mathcal{G}_n(r) \subseteq \mathcal{G}$, the radius of $\mathcal{G}_n(r)$ is upper bounded by $\min\{\sqrt{r}, 1\}$.

Dudley's inequality (Lemma B.1) yields:

$$\mathcal{R}_n(\mathcal{G}_n(r)) \leq \inf_{\alpha \in [0, \min\{\sqrt{r}, 1\}]} \left\{ 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^{\min\{\sqrt{r}, 1\}} \sqrt{\log N(\mathcal{G}; d, \varepsilon/2)} d\varepsilon \right\}. \quad (30)$$

Here, we have used the fact that the inclusion $\mathcal{G}_n(r) \subseteq \mathcal{G}$ implies $N(\mathcal{G}_n(r); d, \varepsilon) \leq N(\mathcal{G}; d, \varepsilon/2)$ by Vershynin [35, Exercise 4.2.10].

Since \mathcal{G} is (B_θ, L_θ, q) -Lipschitz, for any $\varepsilon > 0$, an ε -covering of \mathcal{G} in the d -metric can be constructed from an $\varepsilon/(B_\theta L_\theta)$ -covering of $\mathbb{B}_2^q(1)$ in the Euclidean metric. Therefore, for any $\varepsilon \in (0, 1)$, by the standard volume comparison bound:

$$\begin{aligned} \log N(\mathcal{G}; d, \varepsilon) &\leq \log N\left(\mathbb{B}_2^q(1); \|\cdot\|, \frac{\varepsilon}{B_\theta L_\theta}\right) \\ &\leq q \log\left(1 + \frac{2B_\theta L_\theta}{\varepsilon}\right) \\ &\leq q \log\left(\frac{3B_\theta L_\theta}{\varepsilon}\right). \end{aligned}$$

Putting $R := \min\{\sqrt{r}, 1\}$,

$$\begin{aligned} &\int_0^R \sqrt{\log N(\mathcal{G}; d, \varepsilon/2)} d\varepsilon \\ &\leq \sqrt{q} \left[R\sqrt{\log(6B_\theta L_\theta)} + \int_0^R \sqrt{\log(1/\varepsilon)} d\varepsilon \right] \quad \text{using } \sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \\ &= \sqrt{q} \left[R\sqrt{\log(6B_\theta L_\theta)} + R \int_0^1 \sqrt{\log\left(\frac{1}{R\varepsilon}\right)} d\varepsilon \right] \quad \text{change of variables } \varepsilon \leftarrow \varepsilon/R \\ &\leq \sqrt{q} \left[R\sqrt{\log(6B_\theta L_\theta)} + R\sqrt{\log\left(\frac{1}{R}\right)} + R \right] \quad \text{using } \int_0^1 \sqrt{\log\left(\frac{1}{\varepsilon}\right)} d\varepsilon \leq 1 \\ &\leq R\sqrt{q} \left[\sqrt{\log(6B_\theta L_\theta)} + 2\sqrt{\log\left(\frac{1}{R}\right)} \right] \\ &\leq 2\sqrt{2}R\sqrt{q} \sqrt{\log\left(\frac{6B_\theta L_\theta}{R}\right)} \quad \text{using } \sqrt{a} + \sqrt{b} \leq \sqrt{2}\sqrt{a+b}. \end{aligned}$$

The claim now follows. \square

We complete the proof by upper bounding r_n^* and invoking Theorem B.2. First, observe that by Cauchy-Schwarz, the inequality $\mathbb{E}_n[g^2] \leq \mathbb{E}_n[g]$ for $g \in \mathcal{G}$, and Jensen's inequality:

$$\mathcal{R}_n(\mathcal{H}_n(r)) \leq \sup_{g \in \mathcal{H}_n(r)} \sqrt{\mathbb{E}_n[g^2]} \mathbb{E}_\varepsilon \sqrt{\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2} \leq \sqrt{r}.$$

This bound holds for any $r \geq 0$. Hence, when $r \leq 1/n^2$:

$$\mathcal{R}_n(\mathcal{H}_n(r)) \leq 1/n.$$

On the other hand, when $r > 1/n^2$, by $\mathcal{R}_n(\mathcal{H}_n(r)) \leq \mathcal{R}_n(\mathcal{G}_n(r))$, Proposition B.1, and the inequalities $1/n < \min\{\sqrt{r}, 1\} \leq \sqrt{r}$:

$$\mathcal{R}_n(\mathcal{H}_n(r)) \leq 24\sqrt{2} \sqrt{\frac{q}{n}} \sqrt{r} \sqrt{\log(6B_\theta L_\theta n)}.$$

Hence, the function ϕ_n defined as:

$$\phi_n(r) := \max \left\{ 24\sqrt{2} \sqrt{\frac{q \log(6B_\theta L_\theta n)}{n}} \sqrt{r}, \frac{1}{n} \right\},$$

satisfies $\mathcal{R}_n(\mathcal{H}_n(r)) \leq \phi_n(r)$ for all $r \geq 0$. It is also not hard to see that ϕ_n is a sub-root function (cf. Definition B.1). Therefore, there is a unique solution r_n^* satisfying $\phi_n(r_n^*) = r_n^*$. Now, for any positive constants A, B , the root of $r = \max\{A\sqrt{r}, B\}$ is upper bounded by $\max\{A^2, B\}$. Hence,

$$r_n^* \leq 1152 \frac{q \log(6B_\theta L_\theta n)}{n}.$$

Theorem 4.1 now follows by Theorem B.2. \square

B.4 Proof of Corollary 4.1

Lemma B.2. *Let $B_{\ell,p} := \frac{2}{p+1} \sum_{j=0}^p B_j$ and $L_{\ell,p} := \frac{B_X}{p+1} \sum_{j=0}^p L_j$. Then $B_{\ell,p}^{-1}(\ell_p^{\pi_*} \circ \Pi_{\theta,p})$ is a $(B_\theta, B_{\ell,p}^{-1}L_{\ell,p}, q)$ -Lipschitz parametric function class*

Proof. It suffices to show that

$$\max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_*}(\xi; \pi)\|$$

is $2B_j$ -bounded and $B_X L_j$ Lipschitz with respect to Θ . By definition, we immediately get

$$\begin{aligned} \|\partial_x^j \Delta_t^{\pi_*}(\xi; \pi)\| &= \|\partial_x^j \pi_*(x_t^{\pi_*}(\xi)) - \partial_x^j \pi(x_t^{\pi_*}(\xi))\| \\ &\leq 2 \sup_{\|x\| \leq B_X, \|\theta\| \leq B_\theta} \|\partial_x^j \pi(x, \theta)\| \\ &= 2B_j. \end{aligned}$$

To bound the Lipschitz constant, we iteratively apply the Fundamental Theorem of Line Integrals:

$$\begin{aligned} \partial_x^j \pi(x; \theta_1) - \partial_x^j \pi(x; \theta_2) &= \int_{\theta_2}^{\theta_1} \int_0^x \frac{\partial^{j+2} \pi}{\partial x^{j+1} \partial \theta} (z \otimes \omega) dz d\omega \\ &= \int_{\theta_2}^{\theta_1} \left(\int_0^1 \frac{\partial^{j+2} \pi}{\partial x^{j+1} \partial \theta} (\alpha x \otimes \omega) d\alpha \right) x d\omega \\ &= \left(\int_0^1 \int_0^1 \frac{\partial^{j+2} \pi}{\partial x^{j+1} \partial \theta} (\alpha x \otimes (\theta_2 + \beta(\theta_1 - \theta_2))) d\alpha d\beta \right) x \otimes (\theta_1 - \theta_2). \end{aligned}$$

Taking norms on both sides, we get

$$\begin{aligned} \|\partial_x^j \pi(x; \theta_1) - \partial_x^j \pi(x; \theta_2)\| &\leq \sup_{\|x\| \leq B_X, \|\theta\| \leq B_\theta} \left\| \frac{\partial^{j+2} \pi}{\partial x^{j+1} \partial \theta} \right\| \|x\| \|\theta_1 - \theta_2\| \\ &\leq B_X L_j \|\theta_1 - \theta_2\|, \end{aligned}$$

which establishes that $\|\partial_x^j \Delta_t^{\pi_*}(\xi; \pi)\|$ is $B_X L_j$ -Lipschitz. Recalling that

$$\ell_p^{\pi_*}(\xi; \pi) := \frac{1}{p+1} \sum_{j=0}^p \max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_*}(\xi; \pi)\|,$$

it follows that $\ell_p^{\pi_*}(\xi; \pi)$ is $\frac{2}{p+1} \sum_{j=0}^p B_j$ -bounded and $\frac{B_X}{p+1} \sum_{j=0}^p L_j$ -Lipschitz. \square

Corollary 4.1. *Let the policy class $\Pi_{\theta,p}$ be defined as in (18), and assume that $\pi_* \in \Pi_{\theta,p}$. Let the function class $\ell_p^{\pi_*} \circ \Pi_{\theta,p}$ be defined as in (19), and constants $B_{\ell,p}, L_{\ell,p}$ be defined as above. Let $\hat{\pi}_{\text{TaSIL},p}$ be any empirical risk minimizer (15). Then with probability at least $1 - \delta$ over the initial conditions $\{\xi_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$,*

$$\mathbb{E}_\xi [\ell_p^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},p})] \leq \mathcal{O}(1) B_{\ell,p} \frac{q \log(B_\theta B_{\ell,p}^{-1} L_{\ell,p} n) + \log(1/\delta)}{n}. \quad (20)$$

Proof. This follows by directly applying the constants derived in Lemma B.2 to Theorem 4.1, and using the assumption that $\pi_* \in \Pi_{\theta,p}$ such that $\mathbb{E}_n[\ell_p^{\pi_*}(\cdot; \hat{\pi}_{\text{TaSIL},p})] = 0$. \square

B.5 Proofs of Theorem 4.2 and Theorem 4.3

Before proceeding to the proofs of the main sample complexity bounds, we introduce the following lemma for inverting functions of the form $\log n/n$, adapted from Simchowitz et al. [36, Lemma A.4].

Lemma B.3. *Given $n \in \mathbb{N}$, $n \geq b \log(cn)$ as long as $n \geq 2b \log(2bc)$, where we assume $b, c \geq 1$.*

Proof. We observe by derivatives that $n - b \log(cn)$ is strictly increasing for $n \geq b$. Therefore, it suffices to show $b \log(cn) \leq n$ when $n = 2b \log(2bc)$.

$$\begin{aligned} b \log(2bc \log(2bc)) &= b \log(2 \log(2)bc + 2bc \log(bc)) \\ &\leq b \log((2 \log(2) + 2)(bc)^2) && bc \geq 1 \\ &= 2b \log(\sqrt{2 \log(2) + 2}bc) \\ &< 2b \log(2bc). \end{aligned}$$

□

Theorem B.3 (Full version of Theorem 4.2). *Assume that $\pi_* \in \Pi_{\theta,0}$ and let the assumptions of Theorem 3.1 hold for all $\pi \in \Pi_{\theta,0}$. Let Equation (6) hold with constants $\mu, \alpha > 0$, and assume without loss of generality that $\alpha/\mu \leq 1$, $L_\pi \mu \geq 1/2$. Let $\hat{\pi}_{\text{TaSIL},0}$ be an empirical risk minimizer of $\ell_0^{\pi_*}$ over the policy class $\Pi_{\theta,0}$ for initial conditions $\{\xi_i\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$. Fix a failure probability $\delta \in (0, 1)$, and assume that*

$$n \geq \mathcal{O}(1) \max \left\{ B_{\ell,0} \frac{\kappa_\alpha}{\delta} \log \left(\frac{\kappa_\alpha B_\theta B_{\ell,0}^{-1} L_{\ell,0}}{\delta} \right), B_{\ell,0} \frac{\kappa_\eta}{\delta} \log \left(\frac{\kappa_\eta B_\theta B_{\ell,0}^{-1} L_{\ell,0}}{\delta} \right) \right\},$$

where $\kappa_\alpha := q(\mu/\alpha)^{\frac{1}{1+r}}$, $\kappa_\eta := qL_\pi \mu/\eta$. Then with probability at least $1 - \delta$, the imitation gap evaluated on $\xi \sim \mathcal{D}$ (drawn independently from $\{\xi_i\}_{i=1}^n$) satisfies

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},0}) \leq \mathcal{O}(1) \mu \left(\frac{1}{\delta} \frac{B_{\ell,0} q \log(B_\theta B_{\ell,0}^{-1} L_{\ell,0} n)}{n} \right)^{1+r}.$$

Proof. Applying Corollary 4.1 to the $(B_\theta, B_{\ell,0}, q)$ -Lipschitz parametric function class $B_{\ell,0}^{-1}(\ell_0^{\pi_*} \circ \Pi_{\theta,0})$, we get that with probability at least $1 - \delta/2$ over i.i.d. initial conditions $\xi_i \sim \mathcal{D}^n$,

$$\mathbb{E}_\xi \left[\max_{0 \leq t \leq T-1} \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \right] \leq \mathcal{O}(1) B_{\ell,0} \frac{q \log(B_\theta B_{\ell,0}^{-1} L_{\ell,0} n) + \log(1/\delta)}{n}.$$

Applying Markov's inequality to $\max_t \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|$, for a new draw $\xi \sim \mathcal{D}$, with probability greater than $1 - \delta/2$,

$$\max_{0 \leq t \leq T-1} \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \leq \frac{2}{\delta} \mathbb{E}_\xi \left[\max_{0 \leq t \leq T-1} \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \right].$$

Thus applying a union bound over the two events, we have with probability greater than $1 - \delta$ that

$$\max_{0 \leq t \leq T-1} \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \leq \mathcal{O}(1) B_{\ell,0} q \frac{1}{\delta} \frac{\log(B_\theta B_{\ell,0}^{-1} L_{\ell,0} n) + \log(1/\delta)}{n}, \quad (31)$$

where we absorb numerical constants into $\mathcal{O}(1)$. We want $\max_t \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|$ to satisfy the conditions in (7); that is,

$$\begin{aligned} \max_{0 \leq t \leq T-1} \mu \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|^{1+r} &\leq \alpha, \\ \max_{0 \leq t \leq T-1} 2L_\pi \mu \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|^{1+r} + \|\Delta_t^{\pi_*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| &\leq \eta. \end{aligned}$$

For notational convenience, we further require $\max_t \|\Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \leq 1$, so that

$$\max_t \|\Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|^{1+r} \leq \max_t \|\Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},0})\|.$$

By assumption, since $\alpha/\mu \leq 1$, satisfying the first condition above implies $\max_t \|\Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},0})\| \leq 1$. We observe that for $n \geq \delta^{-1} \log(1/\delta)$ we have $\log n \geq 2 \log(1/\delta)$, thus it suffices to absorb the $\log(1/\delta)$ term into $\log n$. Inserting the generalization bound (31) and shifting n to the right-hand side of the above conditions, we have the following requirements on n :

$$\begin{aligned} n &\geq \mathcal{O}(1) \max \left\{ \left(\frac{\mu}{\alpha} \right)^{1/1+r} B_{\ell,0} q \frac{1}{\delta} \log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n), \right. \\ &\quad \left. \left(\frac{L_{\pi} \mu}{\eta} \right) B_{\ell,0} q \frac{1}{\delta} \log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n) \right\} \\ &=: \mathcal{O}(1) \max \left\{ B_{\ell,0} \kappa_{\alpha} \frac{1}{\delta} \log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n), \right. \\ &\quad \left. B_{\ell,0} \kappa_{\eta} \frac{1}{\delta} \log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n) \right\}, \end{aligned}$$

where we define $\kappa_{\alpha} = q(\mu/\alpha)^{\frac{1}{1+r}}$ and $\kappa_{\eta} = qL_{\pi}\mu/\eta$. Therefore, applying Lemma B.3 on each of the arguments of the maximum, setting $b = B_{\ell,0}\kappa_{\alpha}q/\delta$ (respectively $b = B_{\ell,0}\kappa_{\eta}q/\delta$) and $c = B_{\theta}B_{\ell,0}^{-1}L_{\ell,0}$, we get the following sample complexity bounds. For n satisfying

$$n \geq \mathcal{O}(1) \max \left\{ B_{\ell,0} \frac{\kappa_{\alpha}}{\delta} \log \left(\frac{\kappa_{\alpha} B_{\theta} L_{\ell,0}}{\delta} \right), B_{\ell,0} \frac{\kappa_{\eta}}{\delta} \log \left(\frac{\kappa_{\eta} B_{\theta} L_{\ell,0}}{\delta} \right) \right\},$$

we have with probability greater than $1 - \delta$

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},0}) \leq \mathcal{O}(1) \mu \left(B_{\ell,0} q \frac{1}{\delta} \frac{\log(B_{\theta} B_{\ell,0}^{-1} L_{\ell,0} n)}{n} \right)^{1+r}.$$

This completes the proof. \square

Theorem B.4 (Full version of Theorem 4.3). Assume that $\pi_{\star} \in \Pi_{\theta,p}$, and let the assumptions of Theorem 3.2 hold for all $\pi \in \Pi_{\theta,p}$. Let Equation (10) hold with constants $\mu, \alpha > 0$, and without loss of generality let $(\frac{\alpha}{\mu})^r p! \leq 2$. Let $\hat{\pi}_{\text{TaSIL},p}$ be an empirical risk minimizer of $\ell_p^{\pi_{\star}}$ over the policy class $\Pi_{\theta,p}$ for initial conditions $\{\xi_i\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^n$. Fix a failure probability $\delta \in (0, 1)$, and assume

$$n \geq \mathcal{O}(1) \max_{j \leq p} \max \left\{ B_j \frac{\kappa_{\alpha,j}}{\delta} \log \left(\frac{\kappa_{\alpha,j} B_{\theta} B_j^{-1} B_X L_j}{\delta} \right), B_j \frac{\kappa_{\eta,j}}{\delta} \log \left(\frac{\kappa_{\eta,j} B_{\theta} B_j^{-1} B_X L_j}{\delta} \right) \right\},$$

where $\kappa_{\alpha,j} := (\frac{\mu}{\alpha})^r \frac{pq}{j!}$ and $\kappa_{\eta,j} := (\frac{L_{\theta} p \pi}{(p+1)!} \frac{\mu^{p+1}}{(j!)^{\frac{p+1}{r}}} + \frac{1}{j!}) \frac{pq}{\eta \delta}$. Then with probability at least $1 - \delta$, the imitation gap evaluated on $\xi \sim \mathcal{D}$ (drawn independently from $\{\xi_i\}_{i=1}^n$) satisfies

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},p}) \leq \mathcal{O}(1) \mu \max_{j \leq p} \left(\frac{p}{j! \delta} \frac{B_j q \log(B_{\theta} B_j^{-1} B_X L_j n)}{n} \right)^{1/r}.$$

Proof. Let us first define the following losses on a specific partial:

$$h_j^{\pi_{\star}}(\xi; \pi) := \max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_{\star}}(\xi; \pi)\|.$$

We observe that by definition, $h_j^{\pi_{\star}} \circ \Pi_{\theta,p}$ is $2B_j$ bounded, and $h_j^{\pi_{\star}}$ is $B_X L_j$ -Lipschitz with respect to Θ for $j \leq p$, such that $0.5B_j^{-1}(h_j^{\pi_{\star}} \circ \Pi_{\theta,p})$ is a $(B_{\theta}, 0.5B_j^{-1}B_X L_j, q)$ -Lipschitz loss class. We note that since $\pi_{\star} \in \Pi_{\theta,p}$, we have for any dataset $\{\xi_i\} \subset \mathcal{X}$

$$\mathbb{E}_n [\ell_p^{\pi_{\star}}(\cdot; \hat{\pi}_{\text{TaSIL},p})] =: \frac{1}{p+1} \sum_{j=0}^p \max_{0 \leq t \leq T-1} \|\partial_x^j \Delta_t^{\pi_{\star}}(\xi; \pi)\| = 0,$$

which therefore implies $\mathbb{E}_n[h_j^{\pi^*}(\cdot; \hat{\pi}_{\text{TaSIL},p})] = 0$ for $j \leq p$. We now apply the same proof structure in Theorem 4.2 to each $0.5B_j^{-1}(h_j^{\pi^*} \circ \Pi_{\theta,p})$, where we have with probability greater than $1 - \frac{\delta}{2(p+1)}$ that

$$\mathbb{E}_\xi \left[\max_{0 \leq t \leq T-1} \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},p}) \right\| \right] \leq \mathcal{O}(1) B_j \frac{q \log(B_\theta B_j^{-1} B_X L_j n) + \log\left(\frac{2(p+1)}{\delta}\right)}{n}.$$

Applying Markov's inequality at level $\frac{\delta}{2(p+1)}$, we get with total probability greater than $1 - \frac{\delta}{p+1}$ over a new initial condition $\xi \sim \mathcal{D}$ that $\hat{\pi}_{\text{TaSIL},p}$ satisfies the generalization bound

$$\max_{0 \leq t \leq T-1} \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \hat{\pi}_{\text{TaSIL},p}) \right\| \leq \mathcal{O}(1) B_j \frac{p+1}{\delta} \frac{q \log(B_\theta B_j^{-1} B_X L_j n) + \log\left(\frac{p+1}{\delta}\right)}{n}. \quad (32)$$

For each partial, we want to satisfy the constraints outlined in (11):

$$\begin{aligned} \max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\| \right)^{1/r} &\leq \alpha, \\ \max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\| \right)^{\frac{p+1}{r}} + \frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\| &\leq \eta, \end{aligned} \quad (33)$$

By assumption, we have $\left(\frac{\alpha}{\mu}\right)^r p! \leq 2$, and thus the first condition implies $\max_t \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\| \leq 1$ for all $j \leq p$; in particular, this conveniently ensures $\left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\|^{\frac{p+1}{r}} \leq \left\| \partial_x^j \Delta_t^{\pi^*}(\xi; \pi) \right\|$. Plugging the earlier generalization bound (32) into the above constraints and shifting n to the RHS, and observing like earlier we may absorb the $\log(1/\delta)$ term into the $\log n$ term, we get:

$$\begin{aligned} n &\geq \mathcal{O}(1) \max \left\{ \left(\frac{\mu}{\alpha} \right)^r \frac{1}{j!} B_j \frac{p}{\delta} q \log(B_\theta B_j^{-1} L_j n), \right. \\ &\quad \left. \left(\frac{L_{\partial^p \pi}}{(p+1)!} \frac{\mu^{p+1}}{(j!)^{\frac{p+1}{r}}} + \frac{1}{j!} \right) B_j \frac{p}{\delta} q \log(B_\theta B_j^{-1} L_j n) \right\} \\ &=: \mathcal{O}(1) \max \left\{ B_j \frac{\kappa_{\alpha,j}}{\delta} \log(B_\theta B_j^{-1} L_j n), B_j \frac{\kappa_{\eta,j}}{\delta} q \log(B_\theta B_j^{-1} L_j n) \right\}, \end{aligned}$$

where we define $\kappa_{\alpha,j} = \left(\frac{\mu}{\alpha}\right)^r \frac{pq}{j!}$, $\kappa_{\eta,j} = \left(\frac{L_{\partial^p \pi}}{(p+1)!} \frac{\mu^{p+1}}{(j!)^{\frac{p+1}{r}}} + \frac{1}{j!} \right) \frac{pq}{\eta \delta}$. Therefore applying Lemma B.3, setting $b = B_j \frac{\kappa_{\alpha,j}}{\delta}$ (respectively $b = B_j \frac{\kappa_{\eta,j}}{\delta}$) and $c = B_\theta B_j^{-1} L_j$, for n satisfying:

$$n \geq \mathcal{O}(1) \max \left\{ B_j \frac{\kappa_{\alpha,j}}{\delta} \log\left(\frac{\kappa_{\alpha,j} B_\theta L_j}{\delta}\right), B_j \frac{\kappa_{\eta,j}}{\delta} \log\left(\frac{\kappa_{\eta,j} B_\theta L_j}{\delta}\right) \right\},$$

we have with probability greater than $1 - \frac{\delta}{p+1}$ that the conditions (33) are satisfied. To finish the proof, since we have with probability $1 - \frac{\delta}{p+1}$ that each j th partial difference satisfies the necessary conditions, we union bound over $0 \leq j \leq p$, such that we take a maximum over j for the sample complexity and the resulting imitation gap. This gets us with probability greater than $1 - \delta$, for n satisfying

$$n \geq \mathcal{O}(1) \max_{j \leq p} \max \left\{ B_j \frac{\kappa_{\alpha,j}}{\delta} \log\left(\frac{\kappa_{\alpha,j} B_\theta L_j}{\delta}\right), B_j \frac{\kappa_{\eta,j}}{\delta} \log\left(\frac{\kappa_{\eta,j} B_\theta L_j}{\delta}\right) \right\},$$

that the following bound on the imitation gap holds

$$\Gamma_T(\xi; \hat{\pi}_{\text{TaSIL},p}) \leq \mathcal{O}(1) \max_{j \leq p} \mu \left(\frac{p}{j! \delta} \frac{B_j q \log(B_\theta B_j^{-1} B_X L_j n)}{n} \right)^{1/r}.$$

This completes the proof. \square

C Using finite-differencing to approximate derivatives

C.1 Satisfying Conditions (11) and (12) with approximate derivatives

We recall the closeness conditions on the partials along expert trajectories that guarantee bounds on the imitation gap:

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \right\| \right)^{1/r} \leq \alpha, \quad (11)$$

$$\max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \right\| \right)^{\frac{p+1}{r}} + \frac{2}{j!} \left\| \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \right\| \leq \eta. \quad (12)$$

If we have access to approximate derivatives of the expert $\widehat{\partial_x^j \pi_*}(x)$ such that

$$\left\| \widehat{\partial_x^j \pi_*}(x) - \partial_x^j \pi_*(x) \right\| \leq b < 1$$

for all $x \in \mathbb{R}^d$, then it suffices to tighten the constraints by some function of b such that minimizing with respect to the approximate partial derivatives will still result in the deviation from the true derivatives satisfying the requisite bounds. Let us define

$$\widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) := \partial_x^j \pi(x_t^{\pi_*}(\xi)) - \widehat{\partial_x^j \pi_*}(x_t^{\pi_*}(\xi)),$$

such that

$$\begin{aligned} \left\| \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \right\| &\leq \left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| + \left\| \widehat{\partial_x^j \pi_*}(x_t^{\pi_*}(\xi)) - \partial_x^j \pi_*(x_t^{\pi_*}(\xi)) \right\| \\ &\leq \left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| + b. \end{aligned}$$

Therefore, it suffices to match the approximate partial derivatives such that

$$\begin{aligned} \max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \mu \left(\frac{2}{j!} \left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| \right)^{1/r} &\leq \hat{\alpha}, \\ \max_{0 \leq t \leq T-1} \max_{0 \leq j \leq p} \frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| \right)^{\frac{p+1}{r}} &+ \frac{2}{j!} \left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| \leq \hat{\eta}, \end{aligned}$$

where, provided $\left\| \widehat{\partial_x^j \Delta_t^{\pi_*}}(\xi; \pi) \right\| < 1$:

$$\hat{\alpha} := \left(\alpha^r - \frac{2\mu^r}{j!} b \right)^{1/r}, \quad \hat{\eta} := \eta - \left(\frac{2L_{\partial^p \pi} \mu^{p+1}}{(p+1)!} \left(\frac{2}{j!} \right)^{\frac{p+1}{r}} + \frac{2}{j!} \right) b.$$

A similar bound holds if we also do not have access to the exact derivatives of the learned policy. In practice, these bounds tell us qualitatively that if a sufficiently precise estimate of the derivatives is used, such as through finite differencing, then the imitation gap bounds in Theorem 3.2 still hold.

C.2 Practical approaches for approximating derivatives

Minimizing $\sum_{j=1}^k \left\| \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \right\|$ can be approximated provided π_* can be evaluated at points $\{x_t(\xi) + \delta_i\}_{i=1}^N$ by minimizing the finite difference loss:

$$\ell_{p, \text{FD}}(\xi; \pi, \{\delta_i\}_{i=1}^N) := \max_{1 \leq i \leq N} \left\| \pi_*(x_t(\xi) + \delta_i) - \pi_*(x_t(\xi)) - (\pi(x_t(\xi) + \delta_i) - \pi(x_t(\xi))) \right\|,$$

where the $\{\delta_i\}$ are chosen such that the Taylor expansion

$$\begin{aligned} \sum_{j=1}^p \frac{1}{p!} \partial_x^j \Delta_t^{\pi_*}(\xi; \pi) \cdot \delta_i^{\otimes j} &= \pi_*(x_t(\xi) + \delta_i) - \pi_*(x_t(\xi)) \\ &- (\pi(x_t(\xi) + \delta_i) - \pi(x_t(\xi))) - \frac{R_{p+1}(\delta_i)}{(p+1)!}, \quad \forall 1 \leq i \leq N, \end{aligned}$$

forms a linearly independent system of equations in the derivative parameters. Here, $R_{p+1}(\delta_i)$ denotes the Taylor remainder, which satisfies the inequality $\|R_{p+1}(\delta_i)\| \leq 2L_{\partial^p \pi} \|\delta_i\|^{p+1}$ by Assumption 3.2.

For the case $p = 1$, we can stack the $\{\delta_i\}$ into a matrix S , the finite differences into a matrix M and the remainders into a matrix R to write

$$\partial_x^j \Delta_t^{\pi^*}(\xi; \pi) S = M - R.$$

Provided the $\{\delta_i\}$ are chosen such that S is invertible, the operator norm of $\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)$ can be upper bounded

$$\begin{aligned} \|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| &= \|MS^{-1} - RS^{-1}\| \\ &\leq \|S^{-1}\|(\|M\| + \|R\|) \\ &\leq \|S^{-1}\|(\|M\| + L_{\partial^p \pi} \|S\|^2). \end{aligned}$$

For instance, using a standard basis $S = \varepsilon I$ as the finite difference perturbations yields the following bound on the operator norm:

$$\|\partial_x^j \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{1}{\varepsilon} \|M\| + \varepsilon L_{\partial^2 \pi},$$

where M is the stacked error matrix at the finite differences. Therefore by ensuring sufficiently small ε and finite difference loss, the bound on the Jacobian error can be made arbitrarily small.

Alternatively, if the finite differences δ_i are sampled from a uniform distribution on a sphere of radius ε for each evaluation of $\ell_{1,\text{FD}}$ (i.e, the expert can be cheaply queried during training), Woolfe et al. [37, Theorem 3.15] shows that

$$\|\partial_x^i \Delta_t^{\pi^*}(\xi; \pi)\| \leq \frac{0.8\sqrt{d}}{\zeta^{1/N}} \left(\frac{1}{\varepsilon} \ell_{p,\text{FD}}(\xi; \pi, \{\delta_i\}_{i=1}^N) + \varepsilon L_{\partial^2 \pi} \right),$$

with probability $1 - \zeta$, where d is the dimensionality of the state space. This suggests that provided the expert can be requeried each iteration, $N \ll d$ finite differencing terms can be used.

C.3 Experimental results

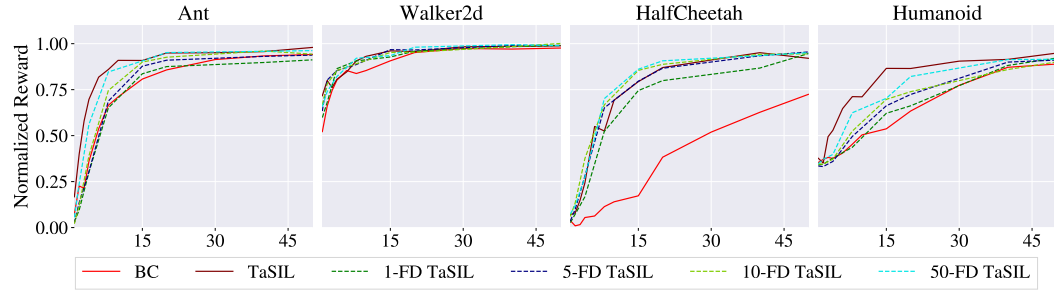


Figure 3: Mean normalized reward for vanilla Behavior Cloning, Behavior Cloning with TaSIL loss, and Behavior Cloning with finite-differencing based TaSIL. The average across 5 random seeds is shown.

We perform several experiments using finite differencing to approximate minimizing the higher order derivatives. Figure 3 shows configurations with 1, 5, 10, and 50 difference vectors across different MuJoCo environments. The different vectors drawn from a uniform distribution over a sphere of radius 0.01. The difference vectors were drawn once for each state-action pair and did not change during training. This was done to simulate the effect of getting progressively closer to using the full standard basis with additional finite differencing terms.

For Walker2d and HalfCheetah with a state dimension of 17, finite difference with a single random vector is sufficient to achieve performance on par with TaSIL using the explicit Jacobians. Humanoid and Ant with higher dimensional observation spaces (376 and 111 dimensions respectively) also show significant improvements the more finite differences are used.

D Additional information for stability experiments

Theorem D.1. For $\eta \in [0, 1)$, the system

$$x_{t+1} = \eta x_t + (1 - \eta) \cdot \frac{\gamma(\|h(x_t) + u_t\|)}{\|h(x_t) + u_t\|} (h(x_t) + u_t), \quad (34)$$

is δ -ISS around $\pi_*(x) = -h(x)$ with class \mathcal{K} function γ .

Proof. We use the shorthand $x_t(\xi_1) := x_t(\xi_1, \{u_k\}_{k=0}^{t-1})$ and $x_t(\xi_2) := x_t(\xi_2, \{0\}_{k=0}^{t-1})$. We can prove this directly using

$$\begin{aligned} & \|x_{t+1}(\xi_1) - x_{t+1}(\xi_2)\| \\ &= \left\| \eta(x_t(\xi_1) - x_t(\xi_2)) + (1 - \eta) \cdot \frac{\gamma(\|h(x_t) + u_t\|)}{\|h(x_t) + u_t\|} (h(x_t) + u_t) \right\| \\ &\leq \eta \|x_t(\xi_1) - x_t(\xi_2)\| + (1 - \eta) \gamma(\|h(x_t) + u_t\|). \end{aligned}$$

Since $x_0(\xi_1) = \xi_1$ and $x_1(\xi_2) = \xi_2$, repeated composition of this upper bound yields

$$\begin{aligned} \|x_t(\xi_1) - x_t(\xi_2)\| &\leq \eta^t \|\xi_1 - \xi_2\| + \sum_{k=0}^{t-1} \eta^{t-1-k} (1 - \eta) \gamma(\|h(x_k) + u_k\|) \\ &\leq \eta^t \|\xi_1 - \xi_2\| + \max_{0 \leq k \leq t-1} \gamma(\|h(x_k) + u_k\|) \\ &= \eta^t \|\xi_1 - \xi_2\| + \gamma \left(\max_{0 \leq k \leq t-1} \|h(x_k) + u_k\| \right). \end{aligned}$$

□

Experiment details The expert MLP has two hidden layers of 32 units each with GELU activations while the learned policy has three hidden layers of 64 units and GELU activations. A tanh nonlinearity was applied to obtain the final policy output. Expert weights were initialized using Lecun Normal initialization LeCun et al. [38] for the kernels and drawn from a normal distribution with $\Sigma = 0.1I$ for the biases. The learned policy weights are initialized using orthogonal initialization for the kernels and zeros for the bias.

For all stability experiments we train on 20 trajectories of length $T = 100$. Initial states were sampled from a standard normal distribution. The state-action pairs are shuffled independently into batches of size 100 and weight updates were performed using the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\varepsilon = 1 \times 10^{-4}$. The training rate was decayed with a cosine learning rate decay using an initial rate of $\alpha = 1 \times 10^{-3}$. We additionally employed ℓ^2 weight regularization with $\lambda = 0.01$. All training is run for 4500 iterations on our internal cluster.

To weight the various derivative terms for the different TaSIL losses we use $\lambda_0 = 1$, $\lambda_1 = 1$, and $\lambda_2 = 10$.

E Additional information for MuJoCo experiments

We use a β -decay-rate of $p = 0.5$ for DAgger and $\alpha = T \text{Tr}[\Sigma_k]$ for DART, the same parameters used by Laskey et al. [7] for their Mujoco experiments. For DART, we use an independent sample of 5 trajectories to update the noise statistics. The same optimization setup from the stability experiments was used, with a batch size of 100, Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 1 \times 10^{-4}$, cosine learning rate scheduling with an initial learning rate of 1×10^{-3} decaying over the entire training duration of 4500 epochs, and ℓ^2 weight regularization with $\lambda = 0.01$.

We train over 4500 epochs for all experiments with a training and test trajectory length of $T = 300$. All TaSIL losses use $\lambda_0 = 1$. $\lambda_1 = 0.01$ is used for the jacobian term in the 1-TaSIL loss.

Similar to Laskey et al. [7], DAgger rollout policies and DART noise statistics were updated sparsely rather than after every trajectory. We performed updates after 1, 5, 20, and 30 trajectories.