

# IMPROVED POLICY OPTIMIZATION FOR ONLINE IMITATION LEARNING

**J. Wilder Lavington**  
University of British Columbia  
jola2372@cs.ubc.ca

**Sharan Vaswani**  
Simon Fraser University  
vaswani.sharan@gmail.com

**Mark Schmidt**  
University of British Columbia  
schmidt@cs.ubc.ca

## ABSTRACT

We consider online imitation learning (OIL), where the task is to find a policy that imitates the behavior of an expert via active interaction with the environment. We aim to bridge the gap between the theory and practice of policy optimization algorithms for OIL by analyzing one of the most popular OIL algorithms, DAGGER. Specifically, if the class of policies is sufficiently expressive to contain the expert policy, we prove that DAGGER achieves constant regret. Unlike previous bounds that require the losses to be strongly-convex, our result only requires the weaker assumption that the losses be strongly-convex with respect to the policy’s sufficient statistics (not its parameterization). In order to ensure convergence for a wider class of policies and losses, we augment DAGGER with an additional regularization term. In particular, we propose a variant of Follow-the-Regularized-Leader (FTRL) and its adaptive variant for OIL and develop a memory-efficient implementation, which matches the memory requirements of FTL. Assuming that the loss functions are smooth and convex with respect to the parameters of the policy, we also prove that FTRL achieves constant regret for any sufficiently expressive policy class, while retaining  $O(\sqrt{T})$  regret in the worst-case. We demonstrate the effectiveness of these algorithms with experiments on synthetic and high-dimensional control tasks.

## 1 INTRODUCTION

Learning to make control decisions online in a stable and efficient manner is important in computer animation (Ling et al., 2020; Zhang & van de Panne, 2018), resource management (Zhou et al., 2011; Ignaciuk & Bartoszewicz, 2010), robotics (Andrychowicz et al., 2020; Xie et al., 2018; Schaal & Atkeson, 2010), and autonomous vehicles (Chen et al., 2020; Sadigh et al., 2016). Online decision making has a variety of challenges: from partial-observability and asymmetric information (Warrington et al., 2021; Choudhury et al., 2018), to function approximation and bootstrapping error (van Hasselt et al., 2018). One common method to avoid some of these problems is through online imitation learning (Ross et al., 2011) (OIL). The OIL setting assumes access to an expert which is known to achieve the desired control objective (e.g. drive safely), and the task is to learn a policy that imitates the behavior of this expert through direct interaction with the environment by the learned policy.

Although there has been substantial progress in practical algorithms for IL such as imitation learning from observations alone (Kidambi et al., 2021; Peng et al., 2018), adversarial IL (AIL) (Ghasemipour et al., 2020; Creswell et al., 2018; Fu et al., 2018), learning from imperfect experts (Sun et al., 2018; Laskey et al., 2017; Sun et al., 2017) or demonstrations (Rengarajan et al., 2022; Reddy et al., 2020; Nair et al., 2018), and learning amortized proposals for planning (Lioutas et al., 2022; Fickinger et al., 2021; Piché et al., 2019), there has been relatively little work on direct policy optimization. Even in areas which sometimes provide guarantees like apprenticeship learning (Shani et al., 2022; Syed & Schapire, 2007; Abbeel & Ng, 2004), and behavioral cloning (BC) (Florence et al., 2021), there is still a large gap between theory and practice.

For OIL, one of the most popular policy optimization algorithms, DAGGER (Ross et al., 2011), minimizes the discrepancy between the learned policy and the expert over all states observed through interaction with the environment. Ross et al. (2011) frame the OIL problem as online convex optimization (OCO) (Hazan, 2019), where the sequence of functions measure the discrepancy between the current policy and the expert. Ross et al. (2011) show that DAGGER is in fact an instance of the follow-the-leader (FTL) algorithm (Hazan et al., 2007) and inherits the FTL guarantees when the discrepancy function is strongly-convex in the parameters of

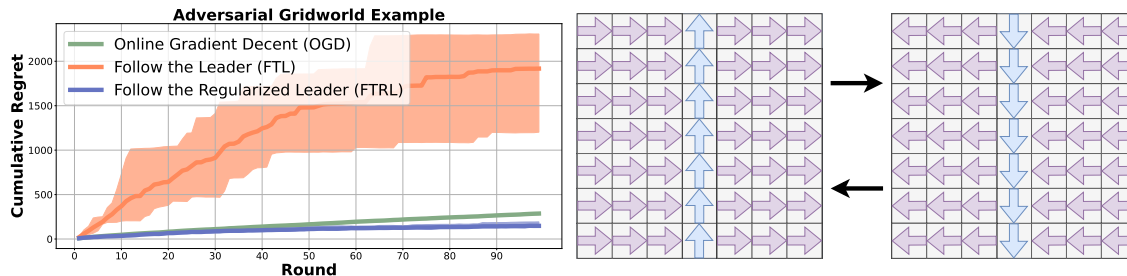


Figure 1: Adversarial Gridworld with a finite-horizon MDP and linear policy parameterization. The agent attempts to match noisy expert feedback in each round: the blue and purple arrows denoting the expert action and the two grids illustrate expert actions in alternate rounds. We observe that FTL has substantially larger cumulative regret, while OGD and FTRL demonstrate better empirical performance (details in Appendix D.1).

the policy. The advantage of the OCO framework is that it actively models adversarial sequences of functions and ensures the resulting algorithms guard against worst case behavior.

However, in OIL the functions are not adversarial. Instead, they are generated by a behavioral policy used to interact with the environment and the function used to measure the discrepancy between the learned and expert policies. Consequently, some OIL algorithms, including DAGGER, have good empirical performance for a broader range of function classes than suggested by the theory. Recent work by Yan et al. (2020) suggests that this theory-practice inconsistency stems from the use of highly expressive policy classes. Assuming that the policy class contains the expert policy, Yan et al. (2020) prove that common OCO algorithms including follow-the-regularized leader (Abernethy et al., 2008), online gradient descent (Zinkevich, 2003), and AdaGrad (Duchi et al., 2011) have better worst-case performance than suggested by the existing theory. However, they (i) only focus on convex functions where modern OIL involves minimizing non-convex loss functions, and (ii) only consider the linearized variants of FTRL such as AdaGrad which are less sample-efficient than their unlinearized counter-parts because they do not take advantage of previous examples. In this work, we address these issues and make the following contributions.

## 1.1 CONTRIBUTIONS

**Follow-the-Leader:** Instead of focusing on the worst-case performance of FTL for arbitrary convex functions, in Section 3, we analyze the theoretical performance of FTL (and thus DAGGER) by exploiting the specific structure in OIL problems. In particular, assuming the class of policies (which we are optimizing over) is sufficiently expressive such that it contains the expert policy, we prove FTL can achieve *constant regret* for OIL problems (Theorem 3.1). Unlike previous work (Yan et al., 2020), this result justifies the superior empirical performance of FTL, and does not require convexity or smoothness with respect to the policy parameterization. Furthermore, we show the use of expressive policy classes can also improve the computational complexity of FTL, making it more robust to hyperparameter tuning. Our analysis shows that much of the empirical success of DAGGER might be due both to its use of specific loss functions as well as expressive policy classes.

However, when the policy class is not rich enough or the agent is not provided with sufficient state information, DAGGER can result in linear regret and poor empirical performance (Warrington et al., 2021; Choudhury et al., 2018). As a simple example, the grid-world in Fig. 1 shows how DAGGER (denoted by its corresponding online optimization algorithm FTL) might exhibit poor oscillatory behavior. This is reminiscent of the counter-example for FTL in the OCO setting in the absence of strong-convexity (Shalev-Shwartz et al., 2012). In the OCO literature, adding a regularization term is the standard way to remedy such oscillatory behavior (Abernethy et al., 2008).

**Follow-the-Regularized-Leader:** Analogously, we use the follow-the-regularized-leader (FTRL) algorithm (Abernethy et al., 2008) in the OIL setting (Section 4). FTRL generalizes FTL, helping it guard against adversarial examples similar to Fig. 1. Unlike FTL, our FTRL analysis assumes loss functions are smooth and convex with respect to the policy parameterization. In Theorem 4.4, we prove that FTRL can obtain constant regret if the policy class is sufficiently expressive and contains the expert policy. In the absence of an expressive policy class, FTRL still results in sublinear regret, improving over FTL in this setting. However, unlike FTL which is parameter-free, FTRL requires that its regularization strength be set according to unknown problem-dependent constants, which can result in poor empirical performance. Consequently,

we make use of the adaptive FTRL (AdaFTRL) algorithm (McMahan, 2017) in the OIL setting. Using a proof technique similar to Vaswani et al. (2020); Yan et al. (2020); Levy et al. (2018); Xie et al. (2020), for smooth, convex loss functions, we prove that AdaFTRL obtains the same regret guarantees as FTRL. We additionally show that unlike FTRL, AdaFTRL does not require the knowledge of problem-dependent constants (Theorem 4.4).

**Experimental Evaluation:** In Section 5, we evaluate the algorithms for both continuous (Todorov et al., 2012) and discrete control (Mnih et al., 2013). Our experiments demonstrate the superior performance of FTL and (Ada-)FTRL, methods that update the policy by utilizing all the past data (so-called “offline” updates). Our experiments also indicate the benefit of functional regularization, with the FTRL variants often outperforming DAGGER (FTL) in terms of either average cumulative loss or average return. An accompanying codebase can be found here: [Improved-Policy-Optimization-for-Online-Imitation-Learning.git](https://github.com/google-research/google-research/tree/master/Improved-Policy-Optimization-for-Online-Imitation-Learning). This code includes all algorithms and baselines discussed in the paper, as well as the additional experiments discussed in Appendix D.

## 2 PROBLEM FORMULATION

We consider an infinite-horizon discounted Markov decision process (MDP) (Bertsekas, 2019; Sutton & Barto, 2018; Puterman, 1994) denoted by  $M$ , and defined by the tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \rho, \gamma \rangle$  where  $\mathcal{S}$  is the set of states,  $\mathcal{A}$  is the action set,  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  is the transition distribution,  $\rho \in \Delta_{\mathcal{S}}$  is the initial state distribution and  $\gamma \in [0, 1)$  is the discount factor. Here,  $\Delta_{\mathcal{S}}$  and  $\Delta_{\mathcal{A}}$  refer to the  $|\mathcal{S}|$ -dimensional and  $|\mathcal{A}|$ -dimensional probability densities respectively. In reinforcement learning (RL), we wish to maximize a reward function denoted by  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . The expected discounted return or *value function* of a policy  $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  is defined as  $V^{\pi}(s) = \mathbb{E}_{a_0, s_1, a_1 \dots} [\sum_{\tau=0}^{\infty} \gamma^{\tau} r(s_{\tau}, a_{\tau}) | s]$ , where  $a_{\tau} \sim \pi(\cdot | s_{\tau})$ , and  $s_{\tau+1} \sim \mathcal{P}(s_{\tau+1} | s_{\tau}, a_{\tau})$  and  $V^{\pi}(\rho) := \mathbb{E}_{s_0 \sim \rho} V^{\pi}(s)$ . A policy  $\pi$  induces a measure  $d^{\pi}$  over states such that  $d^{\pi}(s) = (1 - \gamma) \sum_{\tau=0}^{\infty} \gamma^{\tau} \Pr^{\pi}[s_{\tau} = s | s_0 \sim \rho]$ , where  $\Pr^{\pi}(s_{\tau} = s | s_0 \sim \rho)$  is the visiting probability of  $s$  when playing policy  $\pi$  starting from  $s_0 \sim \rho$ . Given a class of policies  $\Pi$ , the objective is to return a policy that maximizes the value function,  $\max_{\pi \in \Pi} V^{\pi}(\rho)$ .

### 2.1 REDUCTION TO ONLINE CONVEX OPTIMIZATION

Choosing a reward function which is easy to learn from, and that achieves a desired engineering goal, can be difficult (Dulac-Arnold et al., 2021; Hadfield-Menell et al., 2017; Sadigh et al., 2017). Therefore, control engineers often use expert supervision to directly learn an optimal policy. Such experts can make execution of controllers online computationally cheaper at test time Lioutas et al. (2022), or to “warm-start” learning for a complicated control task (Liu et al., 2021). In all cases, given access to such an expert policy  $\pi_e$ , the aim of OIL is to output a policy that imitates the expert. In particular, if the divergence  $D : \Delta_{\mathcal{A}} \times \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$  measures the discrepancy between two policy distributions (for example the KL or Wasserstein divergence),

$$\pi^* = \min_{\pi \in \Pi} \mathbb{E}_{s \sim d^{\pi}} [D(\pi(\cdot | s), \pi_e(\cdot | s))]. \quad (1)$$

That is to say, we search for a policy which minimizes the divergence between expert and agent policies. By imitating the expert, OIL intends to learn a policy that achieves a high return such that  $V^{\pi^*}(\rho) \approx V^{\pi_e}(\rho)$ . Since we cannot compute  $d^{\pi}$  or differentiate through it in general, OIL iteratively samples states from  $d^{\pi}$ , and solves the following optimization problem at iteration  $t \in [T]$  (Ross et al., 2011; Yan et al., 2020),

$$\pi_{t+1} = \arg \min_{\pi \in \Pi} l_t(\pi) = \arg \min_{\pi \in \Pi} \mathbb{E}_{s \sim d^{\pi_t}} [D(\pi(\cdot | s), \pi_e(\cdot | s))]. \quad (2)$$

In this paper, we assume that  $\Pi$  consists of policies that are realizable through a set of sufficient statistics, by a model parameterized with  $w \in \mathcal{W}$ . We use  $\pi_w$  to refer to the parametric realization of  $\pi$ , with the choice of the policy parameterization implicit in the  $\pi_w$  notation. For example, a *linear policy parameterization* often assumes access to features  $\phi_{s,a} \in \mathbb{R}^d$ , and assumes that there exists a  $w \in \mathbb{R}^d$  such that,  $\pi(a | s) \propto \exp(\langle w, \phi_{s,a} \rangle)$ . Notably, if the divergence is convex with respect to the parameterization  $w$ , then each *loss function*  $l_t$  is also convex in  $w$ , and OIL can be recast as an *online convex optimization problem* (Hazan, 2019). Note that here, the functions  $l_t$  are not independent and identically distributed, but instead are generated by a complex interaction between the policy  $\pi_t$  and the MDP at every iteration. In online optimization, we are tasked with finding a sequence of policies parameterized by  $\{w_1, w_2, \dots, w_T\}$  that minimize the sequence of loss functions  $l_t$ . Given  $T$ , the performance of an online optimization algorithm that produces a sequence  $\{w_1, w_2, \dots, w_T\}$

**Algorithm 1** Online optimization for OIL

- 
- 1: **Input:** Policy parameterization  $\pi_w$ , Initial policy  $\pi_{w_0}$      $\triangleright \pi_{w_0}$  has full support over the expert actions.
  - 2: **for**  $t = 1, 2, \dots, T$  **do**
  - 3:    Starting from  $s \sim \rho$ , roll-out policy  $\pi_{w_t}$ .     $\triangleright$  Such interactions can be batched together.
  - 4:    Construct the loss  $l_t(w)$  defined in Eq. (2).
  - 5:    Update  $w_{t+1}$  using an online optimization algorithm.     $\triangleright$  e.g. OGD, AdaGrad, FTRL, FTL.
  - 6: **end for**
  - 7: **Output:** Control policy parameters  $w_T$
- 

is measured in terms of its regret  $R(T)$  defined as:

$$R(T) := \sum_{t=1}^T l_t(w_t) - \min_{w \in \mathcal{W}} \left[ \sum_{t=1}^T l_t(w) \right]. \quad (3)$$

$R(T)$  measures the sub-optimality of the algorithm compared to the best performance in hindsight. We define  $\mathcal{W}_T^* := \arg \min_{w \in \mathcal{W}} \sum_{t=1}^T l_t(w)$  as the best parameter in hindsight. Algorithms that achieve sublinear regret for which  $\lim_{T \rightarrow \infty} R(T)/T = 0$  are referred to as *no-regret algorithms*. It is important to note that common algorithms (Orabona, 2019) achieving sublinear regret do so for any (potentially adversarial) sequence of loss functions. Since it is difficult to correctly account for the interdependence of loss function and policy (which is to say, the policy plays a role in the generation of the next observed loss function), no-regret algorithms instead guarantee performance by safeguarding against worst-case behavior. In the next section, we focus on one such no-regret algorithm, follow-the-leader (FTL), and analyze its performance in OIL.

### 3 FOLLOW-THE-LEADER

In this section, we begin by stating the Follow-the-Leader (FTL) update and then characterize its theoretical performance on OIL problems. The basic FTL update is given by:

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} F_t(w) := \sum_{i=1}^t l_i(w), \quad (4)$$

where  $l_i(w)$  is defined following Eq. (2). This algorithm is desirable because it is parameter free, and makes use of “offline” updates (Schulman et al., 2017; Mnih et al., 2016; Lillicrap et al., 2016; Schulman et al., 2015; Degris et al., 2012) by taking advantage of examples gathered during all previous interactions. These offline updates allow the algorithm to improve the policy without further interactions from the environment, and is important in settings where gathering environment interactions is expensive. In the general OCO framework, when the loss functions  $l_t$  are strongly-convex in  $w$ , FTL achieves  $O(\log(T))$  regret, but will incur  $\Omega(T)$  regret in the absence of strong-convexity (Hazan et al., 2007). However, these results do not capture the empirical success of FTL (e.g. DAGGER) used in conjunction with complex policy parameterizations like neural networks for which the loss functions are non-convex in  $w$  (Warrington et al., 2021).

To make progress towards addressing this discrepancy, we assume that (i) the policy class  $\Pi$  is sufficiently expressive so as to contain the expert policy, (ii) the optimization problem in Eq. (4) can be solved exactly, and (iii) the divergence  $D(\pi, \pi_e)$  has a unique minimizer and is bounded in the sufficient statistics of  $\pi$ . Crucially, *we do not make any assumptions about the policy parameterization*. Assumption (i) is true when using expressive policy classes like neural networks, while assumption (ii) relates to the supervised learning problem in Eq. (4). (ii) can be satisfied if the objective satisfies a gradient domination condition in  $w$ , like for example, the PL inequality (Karimi et al., 2016). Lastly, assumption (iii) is typically true for the divergence-distribution pairs used in practice. For example, consider a continuous state-action space, where for a fixed state  $s$  and  $\pi(a|s) = \mathcal{N}(a; \mu^\top s, I)$  and  $\pi_e(a|s) = \mathcal{N}(a; \mu_e^\top s, I)$ , then  $D(\pi, \pi_e) = 1/2 \|\mu^\top s - \mu_e^\top s\|^2$ . Under these assumptions, we prove (in Appendix B) FTL incurs constant regret.

**Theorem 3.1** (Follow-the-leader - Online Imitation Learning ). *Under the following assumptions: (i) the policy class  $\Pi$  is sufficiently expressive so as to contain the expert policy, the (ii) optimization problem in Eq. (4) can be solved exactly, and (iii) the divergence  $D(\pi, \pi_e)$  has a unique minimizer and is bounded in the sufficient statistics of  $\pi$ , FTL (Eq. (4)) obtains the following regret guarantee in the OIL setting,*

$$R(T) \leq \frac{C}{1-\gamma} \quad \text{where, } C := \max_{\tau, t} \left\{ \mathbb{E}_{s \sim p_t^\tau(s)} [D(\pi(\cdot|s), \pi_e(\cdot|s))] \right\} < \max_s D(\pi(\cdot|s), \pi_e(\cdot|s)).$$

where  $p_t^\tau(s)$  is the probability of reaching state  $s$  at time-step  $\tau$  using policy  $\pi_t$ .

The above result implies that FTL can take advantage of an expressive policy class, and obtain *constant regret* in the OIL setting. Unlike the general online convex optimization results that require strong-convexity in  $w$  and imply a logarithmic regret for FTL (for completeness, we include these proofs in Appendix C.7-C.8), the above theorem doesn't require the strong-convexity of  $l_t(w)$  and is in fact independent of parameterization. Next, we consider the practical implementation of FTL and discuss the advantages of using expressive policy classes in conjunction with FTL updates.

**Solving subproblem in Eq. (4):** If  $\pi_e \in \Pi$ , then there exists a  $w_e$  s.t.  $\pi_e = \pi(w_e)$ . Since  $D(\pi_e, \pi(w_e)) = 0$ , for all  $t$ ,  $l_t(w_e) = 0$ . Hence, the finite-sum problem in Eq. (4) satisfies the *interpolation* (Vaswani et al., 2019a; Ma et al., 2018) property. In this case, stochastic gradient descent (using a randomly sampled  $l_i$ ) matches the convergence rate of deterministic gradient descent on  $F_t$ . For example, if  $F_t(w)$  satisfies the PL property, then it can be minimized to an  $\epsilon$ -error in  $O(\log(1/\epsilon))$  gradient evaluations, making the cost of the FTL update independent of  $t$ . We note that modern machine learning models (e.g. deep neural networks) used for OIL are sufficiently expressive and can ensure that the expert policy is contained in the resulting policy class. Furthermore, note that under interpolation, we can solve the subproblem using SGD with a stochastic line-search (Vaswani et al., 2019b), making the FTL update fully “parameter-free”.

We thus see that having an expressive policy class has a statistical (smaller number of interactions with the environment) as well as a computational advantage (small number of iterations to solve the sub-problem for each update). But in cases where  $\Pi$  does not contain  $\pi_e$ , we have demonstrated (see Fig. 1) that FTL can result in linear regret and poor empirical performance. In order to remedy this issue, we consider a more general class of algorithms known as follow-the-regularized-leader (FTRL) (Abernethy et al., 2008).

## 4 FOLLOW-THE REGULARIZED-LEADER

Recall that for an expert with noisy feedback, FTL can lead to oscillations resulting in large cumulative regret (Fig. 1). We propose to use regularization to stabilize the behavior of FTL. In particular, we analyze the Follow-the-Regularized-Leader (FTRL) algorithm. We first state the FTRL update and then reformulate it for a more scalable practical implementation. For smooth, convex losses, we quantify the regret of FTRL and its adaptive variant in Theorem 4.3 and Theorem 4.4 respectively. The FTRL update (specifically the proximal variant given by Abernethy et al. (2008)) can be defined as:

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} F_t(w) + \psi_t(w) := \left[ \sum_{i=1}^t l_i(w) + \sum_{i=1}^t \frac{\sigma_i}{2} \|w - w_i\|^2 \right]. \quad (5)$$

Note that FTRL can be used in conjunction with other regularizers (Orabona, 2019), but we focus on the squared Euclidean distance throughout this paper. The above update reduces to FTL (Eq. (4)) when  $\sigma_i = 0$  for all  $i$ . Eq. (5). Our analysis uses a proximal regularization term similar to McMahan (2017), though other variants also exist. Note that a naive implementation of Eq. (5) requires storing all the previous parameters  $(w_1, w_2, \dots, w_t)$ . This issue is exacerbated when using large, complex models to parameterize the policy. Using FTRL for continual learning also results in the same problem (Kirkpatrick et al., 2017), and is tackled heuristically. Instead, we reformulate the update in Eq. (5) as follows:

**Proposition 4.1** (Reformulation). *Defining  $\eta_t := 1/(\sum_{i=1}^t \sigma_i)$ , the update in Eq. (5) can be reformulated (proof in Appendix C) as:*

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} \left[ \sum_{i=1}^t l_i(w) - \langle w, \sum_{i=1}^{t-1} \nabla l_i(w_t) \rangle + \frac{1}{2\eta_t} \|w - w_t\|^2 \right]. \quad (6)$$



Unlike Eq. (5), this update does not have a memory requirement which increases with the number of iterations and model size. That is, if  $m$  is the model-size, then Eq. (5) requires  $O(mT)$  memory, while the reformulated update can be implemented using only  $O(m + T)$  memory (same as FTL). We note such reformulations are not unique, and choosing one reformulation over another could lead to drastically different solutions in settings where the inner optimization problem defined by equation 5 is non-convex, or solved inexactly. While we leave a theoretical discussion on this topic to future work, we include an additional reformulation for comparison which we refer to as Alt-FTRL,

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} \left[ \sum_{i=1}^t l_i(w) + \frac{1}{2\eta_t} \|w\|^2 - w^\top \left[ \sum_{i=1}^{t-1} w_t \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right] \right] \right]. \quad (7)$$

This reformulation averages over the previous parameters instead of their gradients, and is included in the empirical comparison in Section 5. Next, we describe how interpolation improves the computational efficiency of FTRL.

**Solving subproblem in Eq. (5):** Similar to Eq. (4), if  $\pi_e \in \Pi$ , the finite-sum in  $F_t(w)$  satisfies the interpolation property defined in Section 3. In this case, proximal stochastic gradient descent (using a randomly sampled  $l_i$  and a proximal operator with  $\psi_t(w)$  discussed in Appendix C) matches the convergence rate of deterministic gradient descent (Cevher & Vü, 2019). Hence, similar to Eq. (4), Eq. (5) can be solved efficiently. Unlike the result in Section 3, here, we assume that the loss functions are  $L$ -smooth and convex in  $w$  (see Appendix A for formal definitions and Appendix C.6 for proofs under the convex, non-smooth but Lipschitz setting). The subsequent results also assume that  $\mathcal{W}$  is a convex compact set of diameter  $D$ , meaning  $\sup_{x,y \in X} \|x - y\| \leq D$ . We use the following definition to quantify the degree to which the interpolation property is satisfied.

**Definition 4.2** (Interpolation-Error). For a fixed iteration  $t$ , if  $w_t^* := \arg \min_{w \in \mathcal{W}} l_t(w)$ , and for a  $w^* \in \mathcal{W}^*$ ,

$$\epsilon_t^2 := \min_{w^* \in \mathcal{W}^*} \{l_t(w^*) - l_t(w_t^*)\}. \quad (8)$$

If  $\pi_e \in \Pi$ , then  $w_e \in \mathcal{W}^*$ . We know that  $l_t(w_e) = 0$  for all  $t$ , and since each  $l_t$  is lower-bounded by zero, it implies that  $\epsilon_t^2 = 0$  for all  $t$ . Hence,  $\epsilon_t^2$  is a measure of the expressivity of the policy class that is induced by the  $w \rightarrow \pi$  mapping and the set  $\mathcal{W}$ . In the following theorem (proved in Appendix C), under the above assumptions, we show that FTRL incurs sublinear regret regardless of whether  $\pi_e \in \Pi$ .

**Theorem 4.3** (FTRL - Smooth + Convex). *Assuming each  $l_t$  is (i)  $L$ -smooth, (ii) convex, FTRL (Eq. (6)) for  $\eta_t = \min \left\{ (\sum_{t=1}^T \epsilon_t^2)^{-1/2}, \frac{1}{2L} \right\}$  for all  $t$ , achieves the following regret*

$$R(T) \leq \sum_{t=1}^T \left[ \frac{\eta_t \|\nabla l_t(w_t)\|^2}{2} \right] + \frac{D^2}{2\eta_T} \leq 2D^2L + (D^2 + 2L) \sqrt{\sum_{t=1}^T \epsilon_t^2}.$$

This result follows similar proof techniques established in works like Orabona (2019); Ghadimi & Lan (2013), however unlike previous results for FTRL, we present a bound which explicitly accounts for the level of interpolation similar to (Loizou et al., 2021; Vaswani et al., 2020). If  $\pi_e \in \Pi$ ,  $\epsilon_t = 0$  for all  $t$ , and FTRL achieves constant regret similar to FTL. For non-zero  $\epsilon_t$ , FTRL still incurs sublinear  $O(\sqrt{T})$  regret. However, we note that unlike Theorem 3.1, the above result requires the loss functions to be convex and smooth. Unfortunately, the above result requires setting  $\eta$  according to  $L$  and  $\epsilon_t$ , both of which are typically unknown in practice. To address this issue, we use an adaptive variant of FTRL (Joulani et al., 2020; McMahan, 2017) and characterize its regret bound in the following theorem (proved in Appendix C).

**Theorem 4.4** (AdaFTRL - Smooth + Convex). *Assuming each  $l_t$  is (i)  $L$ -smooth, (ii) convex, FTRL (Eq. (6)) for  $\eta_t = \alpha / \sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}$ , achieves the following regret*

$$R(T) \leq \sum_{t=1}^T \left[ \frac{\alpha \|\nabla l_t(w_t)\|^2}{2 \sum_{i=1}^t \|\nabla l_i(w_i)\|^2} \right] + \frac{D^2}{2\alpha} \sqrt{\sum_{t=1}^T \|\nabla l_t(w_t)\|^2} \leq 2L \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right)^2 + \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T \epsilon_t^2}.$$

Observe that the above bound holds for any finite value of  $\alpha$ , though the upper-bound is minimized when  $\alpha = D$ . The only difference with the FTRL update used in Theorem 4.3 is the choice of  $\eta_t$ . Hence, we can continue to use the reformulated update in Proposition 4.1. Furthermore, in smooth convex settings, AdaFTRL achieves the same regret as in Theorem 4.3 without the knowledge of problem-dependent constants like  $L$  or  $\epsilon_t$ . The proof of Theorem 4.4 is similar to AdaGrad (Vaswani et al., 2020; Levy et al., 2018). We conclude by showing that FTRL (and AdaFTRL) generalizes OGD (and AdaGrad) respectively. In particular, if we were to use the linearized losses, meaning that  $l_s(w) = \langle \nabla l_s(w_s), w \rangle$  in Proposition 4.1, then, by definition of  $w_{t+1}$ , setting the gradient to zero,  $\sum_{i=1}^t \nabla l_i(w_i) - \sum_{i=1}^{t-1} \nabla l_i(w_i) + \eta_t(w_{t+1} - w_t) = 0$ ,  $\implies w_{t+1} = w_t - \eta_t \nabla l_t(w_t)$ , which recovers the OGD (and AdaGrad) update with the corresponding choice of  $\eta_t$ .

## 5 EXPERIMENTS

In this section, we compare FTL, FTRL, Alt-FTRL, AdaFTRL, online gradient descent (Zinkevich, 2003) (OGD) and AdaGrad (Duchi et al., 2011) in terms of both the average cumulative loss equal to  $1/t \sum_{s=1}^t l_s(w_s)$  and the policy return  $V_t^\pi(\rho)$ . Every round consists of  $M$  interactions with the environment and we evaluate each algorithm for a total of 25000 (10000 for Atari) environment interactions, meaning that  $T = 25000/M$ . Throughout the main paper, we use  $M = 1000$ , and defer the results for  $M = 100$  to Appendix D). In order to tune the "outer-learning-rate"  $\eta$  for OGD, AdaGrad, and the FTRL variants, we do a grid-search over  $\eta \in [10^{-5}, 10^{-4}, \dots, 10^5]$  on one of the environments, and take the top three step-sizes that has the minimum average cumulative loss over the course of 2000 environment interactions under  $M = 100$ . We then evaluate these three step-sizes over 25000 interactions, and take the best of the three in terms of average cumulative loss. For the off-policy methods (FTL and the FTRL variants), we also search over the space of "inner-learning-rates"  $\alpha \in [10^{-5}, 10^{-4}, \dots, 10^5]$ , and run two optimization procedures, Adam (Kingma & Ba, 2014), and SLS (Loizou et al., 2021). For FTRL and Alt-FTRL, we set  $\eta_t = \alpha/\sqrt{t}$ . For FTL, FTRL and its variants, we used gradient descent to solve the subproblems for each update. For solving each subproblem, we used a maximum of 1000 iterations terminating the optimization when the gradient norm was sufficiently small ( $10^{-8}$ ).

### 5.1 CONTINUOUS CONTROL ON MUJOCO

We evaluate the algorithms for continuous control tasks (with continuous state and action spaces) in Mujoco suite (Todorov et al., 2012), and build and train models using pytorch (Paszke et al., 2019). In particular, we consider the Hopper and Walker-2D environments where the task is to learn a policy that can imitate the expert policy trained using reinforcement learning. The expert policy uses a neural network parameterization and is trained using soft actor-critic (Haarnoja et al., 2018). For each environment, we report the performance of each method when the loss function  $l_t(w)$  is either the  $l_2$  or  $l_1$  loss. All results are averaged over 3 runs, and we report the mean and relevant quantiles. The policy corresponds to a multivariate Gaussian distribution with a fixed diagonal covariance and the mean parameterized by either a linear or neural network model. The neural network architecture is the same as that of the expert, meaning that in this case, the policy class is sufficiently expressive to include the expert policy. For the linear model, the resulting loss functions are convex, whereas using a neural network parameterization results in non-convex loss functions. Because of its poor empirical performance, we do not plot the OGD in the main paper and defer these plots to Appendix D.

For both parameterizations (Fig. 2 and Fig. 3), we observe that (i) FTL, FTRL, Alt-FTRL, and AdaFTRL consistently outperform OGD and AdaGrad, (ii) FTRL and its variants consistently outperform FTL in terms of the average cumulative loss, (iii) AdaFTRL improves over both FTRL and Alt-FTRL in terms of average cumulative loss, (iv) FTL has good performance for the non-strongly-convex  $l_1$  loss in the linear case where interpolation is not necessarily satisfied, and (v) good performance with respect to the average loss metric does not imply good return (for example, with the linear model and L2 loss, AdaGrad matches the average loss of the other methods, but has poor performance with respect to the cumulative return). We conclude that (i) "offline" updates used in FTL, FTRL and its variants result in superior empirical performance and (ii) regularization helps improve the empirical performance with FTRL outperforming FTL in terms of average cumulative loss but not always average cumulative reward, and (iii) FTL performs better compared to what is suggested by the theory.

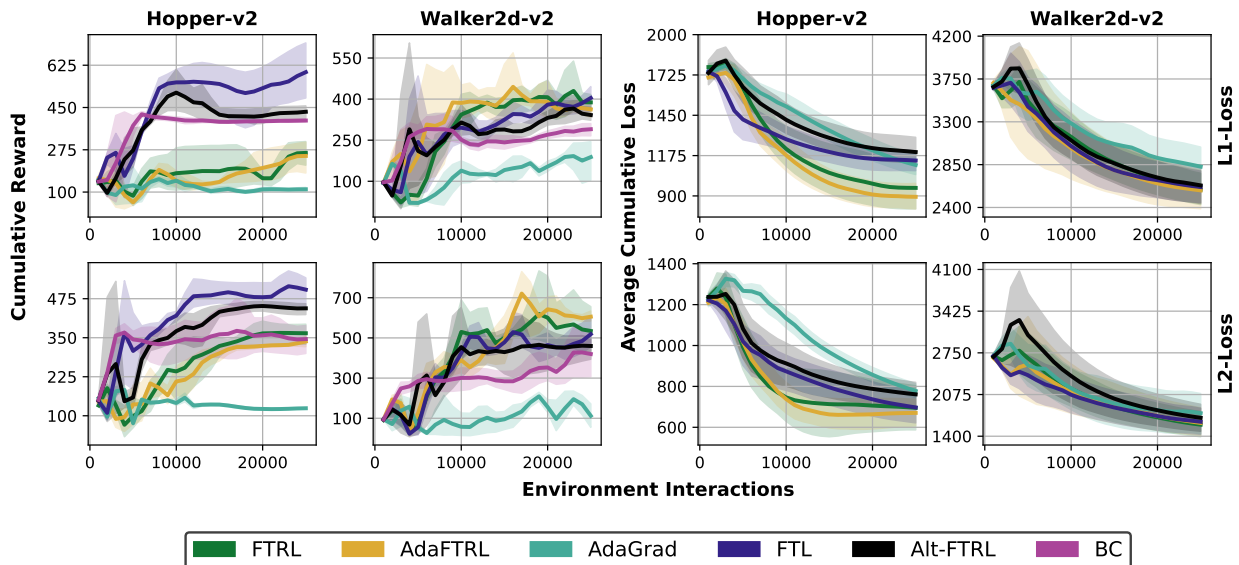


Figure 2: Comparing Adagrad, FTL, FTRL, Alt-FTRL, and AdaFTRL, the four plots on the left display the expected cumulative reward while average cumulative loss displayed on the right. A behavioral cloning baseline (BC) is included to verify the effect of environment interaction under the learned policy. Each line describes the mean, 5% and 95% quantiles as computed by Harris et al. (2020). The plots for the average cumulative loss show that for **linear models**, (i) FTRL, Alt-FTRL and AdaFTRL maintain performance that is as good or better than FTL in terms of average cumulative loss, indicating the benefits of regularization, and (ii) FTL, FTRL, and AdaFTRL do significantly better than AdaGrad, indicating that making “offline” updates are crucial in maintaining performance in online imitation learning.

## 5.2 DISCRETE CONTROL ON ATARI

We evaluate the algorithms for discrete control tasks (with discrete state-action spaces) in the Atari suite (Mnih et al., 2013). In particular, we consider the Pong and Breakout game environments where the task is to learn a policy that can imitate the expert policy trained using reinforcement learning. The expert policy uses a neural network parameterization and is trained using proximal policy optimization (Hill et al., 2018; Schulman et al., 2017). The learned policy corresponds to a categorical distribution parameterized by either a linear model that uses the fixed pretrained features from the reinforcement learning algorithm, or the same neural network architecture as the expert and is learned in an end-to-end fashion. For the linear model (which uses the pretrained feature extractor from the expert), the resulting loss functions are convex, while the end-to-end setup tests the non-convex setting. Because of its poor empirical performance, we again do not plot OGD in the main paper and defer these plots to Appendix D. Notably both settings the policy class is sufficiently expressive so as to include the expert policy.

In Fig. 4, we again observe that for both policy parameterizations, (i) FTL, FTRL and AdaFTRL consistently outperform OGD and AdaGrad, (ii) FTRL and AdaFTRL often dominate FTL in terms of the average cumulative loss, (iii) AdaFTRL has similar performance as its non-adaptive variants in terms of the average cumulative loss, (iv) FTL has good performance for the non-strongly-convex cross-entropy loss, and (v) similar to Section 5.1, good performance with respect to the average loss metric does not imply good return. We again conclude that (i) “offline” updates used in FTL, FTRL and its variants result in superior empirical performance and (ii) regularization helps improve the empirical performance with FTRL outperforming (in average cumulative loss) FTL in the end-to-end setting, and (iii) FTL performs better compared to that suggested by the theory.

## 6 OTHER RELATED WORK

We briefly discussed the most relevant related work in Section 1. We now clarify how OIL relates to other common settings, focusing specifically on imitation learning from observation alone or ILOA (Yan et al., 2020;



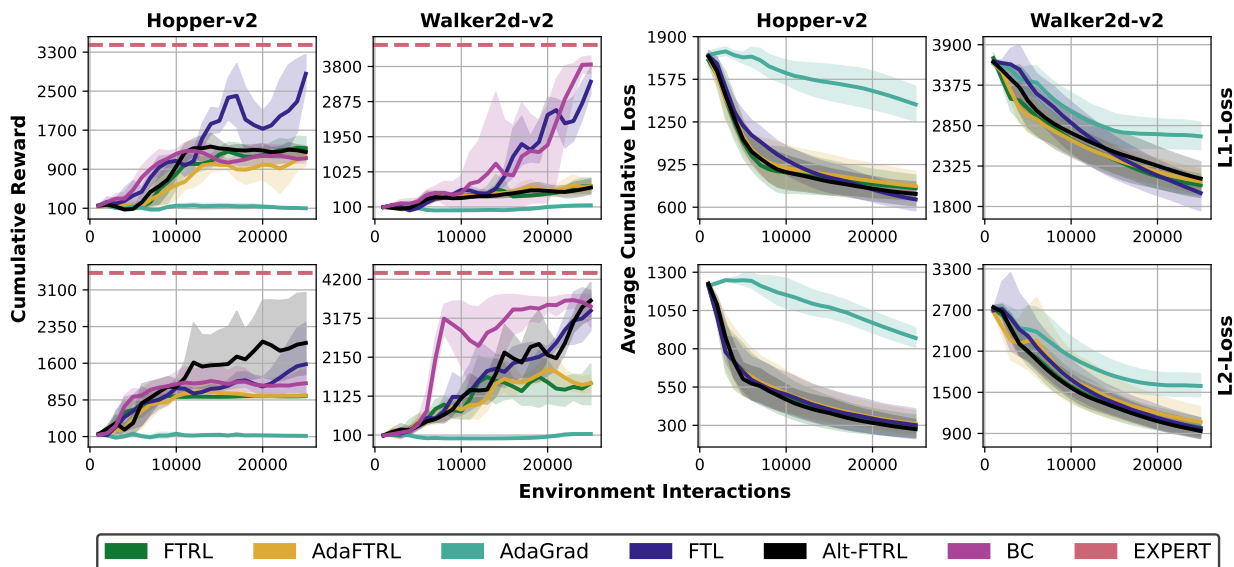


Figure 3: Comparing Adagrad, FTL, FTRL, Alt-FTRL and AdaFTRL, the four plots on the left display the expected cumulative reward in left 4 plots and Average cumulative loss on the right four plots. A behavioral cloning baseline (BC) is included to verify the effect of environment interaction under the learned policy, and expert performance is indicated with a dotted line. Each line presents the mean, 5% and 95% quantiles as computed by Harris et al. (2020). These plots show that for **neural network models**, (i) FTRL, Alt-FTRL, and AdaFTRL maintain performance in terms of cumulative loss that is as good as FTL, even under a complex model class, and (ii) FTL, FTRL, and AdaFTRL all do significantly better than their on-policy counterparts, indicating that making “offline” updates can be crucial in maintaining performance.

Sun et al., 2019). Unlike OIL, ILOA often requires solving difficult sub-problems at every round of interaction. One of the most popular variants of ILOA called *apprenticeship learning - AL* (Shani et al., 2022; Zahavy et al., 2020; Abbeel & Ng, 2004), uses the state-occupancy generated by the expert to construct a reward surface, and solves for the optimal policy under this surface. For example, Syed & Schapire (2007) assume (a) access to expert trajectories and that (b) the underlying unknown reward function is a non-negative linear combination of known “state-features”, and solve a saddle-point problem by alternating between 1) constructing a prospective reward under which the expert occupancy is optimal and the agent’s occupancy is maximally sub-optimal, and 2) using the prospective reward, solve an RL problem via policy optimization. In contrast, the basic OIL setting (Ross et al., 2011; Yan et al., 2020) make no assumption about the unknown reward surface, and instead assumes access to an expert oracle which provides the optimal action given a state. Here, the OIL problem is typically framed as an online optimization problem, can be solved by standard OCO techniques, and notably doesn’t require solving RL problems as a subroutine.

A variety of authors improve upon the framework given by (Syed & Schapire, 2007; Abbeel & Ng, 2004), often taking advantage of theoretical advances in constrained optimization. For example, Frank-Wolfe updates (Zahavy et al., 2020; Abernethy & Wang, 2017), can be used to solve classical variants of the apprenticeship learning problem, and results which extend this framework have even proven some convergence results in non-stationary tabular MDPs (Geist et al., 2022; Zahavy et al., 2021; Zhang et al., 2020). In contrast to these works, our results in Section 4 can handle more complex policy parameterizations (e.g. linear) given that the corresponding losses satisfy the appropriate convexity assumptions, while the FTL results in Section 3 can handle general policy parameterizations under an interpolation assumption (extending similar work in the bandits setting by Degenne et al. (2018)). In some cases by framing the saddle point problem as online mirror-decent (Shani et al., 2022), one can convert the computationally costly saddle-point problem into a more tractable iterative algorithm similar to *adversarial imitation learning* (AIL) (Creswell et al., 2018; Ghasemipour et al., 2020; Fu et al., 2018). Again however, these algorithms require stricter assumptions on the class of MDPs considered, or have no guarantees at all. From a sample-complexity perspective, Baram et al. (2017) shows by assuming access to an expert oracle (like in OIL), the number of environment interactions required to match the expert performance  $\approx 10^5$ , while AIL-like algorithms require up to  $10^6$  even with a similar number of expert examples. Addressing this gap in statistical efficiency between

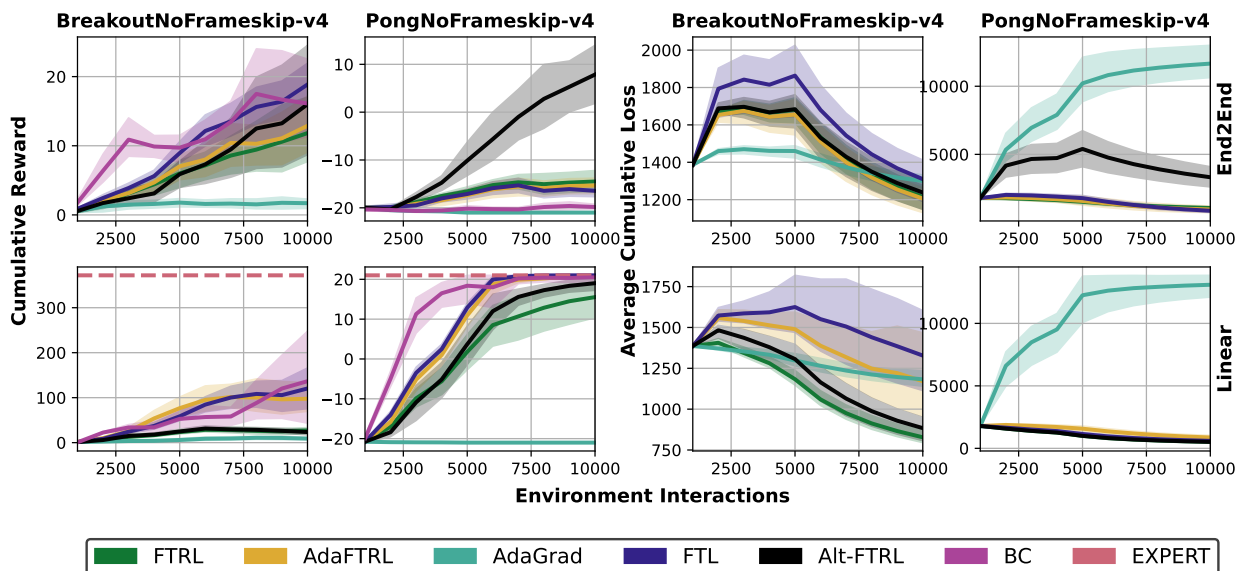


Figure 4: Comparing Adagrad, FTL, FTRL, Alt-FTRL, and AdaFTRL, we plot the the expected cumulative reward on the left 4 plots and the average cumulative loss on the right. Each line presents the mean, 5% and 95% quantiles as computed by [Harris et al. \(2020\)](#). We observe that for the end-to-end setting, (i) FTRL, Alt-FTRL, and AdaFTRL outperform FTL in terms of average cumulative loss, and (ii) FTL, FTRL, Alt-FTRL, and AdaFTRL all do significantly better than AdaGrad, indicating that “offline” updates can be crucial for good empirical performance. A behavioral cloning baseline (BC) is included to verify the effect of environment interaction under the policy, and expert performance is indicated with a dotted line. We also note that the “linear policy class” represents a linear transformation of the experts pretrained feature encoder. This ensures that the learned policy satisfies the interpolation convexity assumptions.

these algorithm classes, and extending results based non-stationary MDPs to continuous state-action spaces represent interesting areas of future research.

## 7 DISCUSSION

We show that in OIL settings (i) algorithms which make use of offline updates (FTL, FTRL, AltFTRL, AdaFTRL) perform better than algorithms which do not (OGD, AdaGrad) and (ii) including regularization can lead to empirical improvements in terms of reward and average cumulative loss. Furthermore, we improved the theoretical results for both FTL and FTRL, showing that both algorithms can achieve constant regret when the policy class is sufficiently expressive and contains the expert policy. Importantly, our guarantees for FTL only require that the losses be strongly-convex with respect to the policy’s sufficient statistics (not its parameters). Our research leaves a host of open questions – (i) does FTRL converge for the IL setting independent of parameterization in a similar fashion to FTL and (ii) can these results be generalized to the standard online learning setting without leveraging OIL structure, and (iii) how do inexact optimization, stochastic gradients, and limited memory affect performance of FTL and FTRL.

## 8 ACKNOWLEDGEMENTS

We thank the anonymous reviewers, whose comments helped improve the clarity of the manuscript. We thank Frederik Kunstner for his invaluable discussion, for providing comments on the manuscript and earlier versions of this work, and for suggesting related material. We also thank Victor Sanches Portella and Betty Shea for additional discussions we believe improved the overall clarity of the work. This research was partially supported by the Canada CIFAR AI Chair Program, the Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grants RGPIN-2015-06068 and the NSERC Post-graduate Scholarships-Doctoral Fellowship 545847-2020.

## REFERENCES

- Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning (ICML)*, Proceedings of Machine Learning Research, pp. 1. PMLR, 2004. URL <https://ai.stanford.edu/~ang/papers/icml04-apprentice.pdf>.
- Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *21st Annual Conference on Learning Theory (COLT)*, pp. 263–274, 2008.
- Jacob D Abernethy and Jun-Kun Wang. On frank-wolfe and equilibrium computation. In *Advances in Neural Information Processing Systems*, 2017. URL <https://papers.nips.cc/paper/2017/hash/7371364b3d72ac9a3ed8638e6f0be2c9-Abstract.html>.
- OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, pp. 3–20, 2020.
- Larry Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of mathematics*, 1966.
- Nir Baram, Oron Anschel, Itai Caspi, and Shie Mannor. End-to-end differentiable adversarial imitation learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 390–399. PMLR, 2017. URL <https://proceedings.mlr.press/v70/baram17a.html>.
- Dimitri P Bertsekas. *Reinforcement learning and optimal control*. Athena Scientific Belmont, MA, 2019.
- Volkan Cevher and Bang Công Vũ. On the linear convergence of the stochastic gradient method with constant step-size. *Optimization Letters*, pp. 1177–1187, 2019.
- Dian Chen, Brady Zhou, Vladlen Koltun, and Philipp Krähenbühl. Learning by cheating. In *Proceedings of the 4th Annual Conference on Robot Learning (CoRL)*, Proceedings of Machine Learning Research, pp. 66–75. PMLR, 2020.
- Sanjiban Choudhury, Mohak Bhardwaj, Sankalp Arora, Ashish Kapoor, Gireeja Ranade, Sebastian Scherer, and Debadepta Dey. Data-driven planning via imitation learning. *The International Journal of Robotics Research*, pp. 1632–1672, 2018.
- Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. Generative adversarial networks: An overview. In *2018 IEEE Signal Processing Magazine*, pp. 53–65. IEEE, 2018.
- Rémy Degenne, Evrard Garcelon, and Vianney Perchet. Bandits with side observations: Bounded vs. logarithmic regret. In *Uncertainty in Artificial Intelligence (UAI)*, 2018. URL <http://auai.org/uai2018/proceedings/papers/182.pdf>.
- Thomas Degris, Martha White, and Richard S. Sutton. Linear off-policy actor-critic. In *International conference on machine learning (ICML)*, Proceedings of Machine Learning Research, pp. 457–464. PMLR, 2012. URL <https://proceedings.neurips.cc/paper/2019/file/0e095e054ee94774d6a496099eb1cf6a-Paper.pdf>.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 2011.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J. Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: Definitions, benchmarks and analysis. *Machine Learning*, pp. 2419–2468, 2021. URL <https://doi.org/10.1007/s10994-021-05961-4>.
- Arnaud Fickinger, Hengyuan Hu, Brandon Amos, Stuart Russell, and Noam Brown. Scalable online planning via reinforcement learning fine-tuning. In *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=D0xGh031I9m>.

- Pete Florence, Corey Lynch, Andy Zeng, Oscar A. Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In *Proceedings of the 5th Annual Conference on Robot Learning (CoRL)*, Proceedings of Machine Learning Research, pp. 158–168. PMLR, 2021. URL <https://proceedings.mlr.press/v164/florence22a.html>.
- Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://openreview.net/forum?id=rkHywl-A->.
- Matthieu Geist, Julien Pérolat, Mathieu Laurière, Romuald Elie, Sarah Perrin, Oliver Bachem, Rémi Munos, and Olivier Pietquin. Concave utility reinforcement learning: The mean-field game viewpoint. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pp. 489–497. International Foundation for Autonomous Agents and Multiagent Systems, 2022. URL <https://arxiv.org/abs/2106.03787>.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, pp. 2341–2368, 2013.
- Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 1259–1277. PMLR, 2020. URL <https://proceedings.mlr.press/v100/ghasemipour20a.html>.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications, 2018. URL <https://arxiv.org/abs/1812.05905>.
- Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J Russell, and Anca Dragan. Inverse reward design. In *Advances in Neural Information Processing Systems*, 2017. URL <https://proceedings.neurips.cc/paper/2017/file/32fdab6559cdfa4f167f8c31b9199643-Paper.pdf>.
- Charles R. Harris, K. Jarrod Millman, Stéfan J van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Pícus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585:357–362, 2020. doi: 10.1038/s41586-020-2649-2.
- Elad Hazan. Introduction to online convex optimization, 2019. URL <https://arxiv.org/abs/1909.05207>.
- Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, pp. 169–192, 2007.
- Ashley Hill, Antonin Raffin, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, Rene Traore, Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, and Yuhuai Wu. Stable baselines. <https://github.com/hill-a/stable-baselines>, 2018.
- Przemysław Ignaciuk and Andrzej Bartoszewicz. Linear–quadratic optimal control strategy for periodic-review inventory systems. *Automatica*, pp. 1982–1993, 2010. URL <https://www.sciencedirect.com/science/article/pii/S0005109810003924>.
- Pooria Joulani, András György, and Csaba Szepesvári. A modular analysis of adaptive (non-) convex optimization: Optimism, composite objectives, variance reduction, and variational bounds. *Theoretical Computer Science*, pp. 108–138, 2020.
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811, 2016.
- Rahul Kidambi, Jonathan Daniel Chang, and Wen Sun. Mobile: Model-based imitation learning from observation alone. In *Advances in Neural Information Processing Systems*, 2021. URL [https://openreview.net/forum?id=\\_Rtm4rYnIIL](https://openreview.net/forum?id=_Rtm4rYnIIL).

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014. URL <https://arxiv.org/abs/1412.6980>.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, pp. 3521–3526, 2017.
- 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 (CoRL)*, Proceedings of Machine Learning Research, pp. 143–156. PMLR, 2017. URL <https://proceedings.mlr.press/v78/laskey17a.html>.
- Kfir Y Levy, Alp Yurtsever, and Volkan Cevher. Online adaptive methods, universality and acceleration. In *Advances in Neural Information Processing Systems*, 2018. URL <https://proceedings.neurips.cc/paper/2018/file/b0169350cd35566c47ba83c6ec1d6f82-Paper.pdf>.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2016. URL <http://arxiv.org/abs/1509.02971>.
- Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. Character controllers using motion vaes. *ACM Transactions on Graphics*, 2020. URL <http://dx.doi.org/10.1145/3386569.3392422>.
- Vasileios Lioutas, Jonathan Wilder Lavington, Justice Sefas, Matthew Niedoba, Yunpeng Liu, Berend Zwartsenberg, Setareh Dabiri, Frank Wood, and Adam Scibior. Critic sequential monte carlo, 2022. URL <https://arxiv.org/abs/2205.15460>.
- Haochen Liu, Zhiyu Huang, Jingda Wu, and Chen Lv. Improved deep reinforcement learning with expert demonstrations for urban autonomous driving, 2021. URL <https://arxiv.org/abs/2102.09243>.
- Nicolas Loizou, Sharan Vaswani, Issam Hadj Laradji, and Simon Lacoste-Julien. Stochastic polyak step-size for sgd: An adaptive learning rate for fast convergence. In *Twenty Fourth International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1306–1314. PMLR, 2021. URL <https://arxiv.org/pdf/2002.10542.pdf>.
- Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the effectiveness of SGD in modern over-parametrized learning. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 3325–3334. PMLR, 2018. URL <https://proceedings.mlr.press/v80/ma18a.html>.
- H Brendan McMahan. A survey of algorithms and analysis for adaptive online learning. *The Journal of Machine Learning Research*, 2017.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013. URL <https://arxiv.org/abs/1312.5602>.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proceedings of The 33rd International Conference on Machine Learning*, Proceedings of Machine Learning Research, pp. 1928–1937. PMLR, 2016. URL <https://proceedings.mlr.press/v48/mniha16.html>.
- Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Overcoming exploration in reinforcement learning with demonstrations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6292–6299. IEEE, 2018.
- Francesco Orabona. A modern introduction to online learning, 2019. URL <https://arxiv.org/abs/1912.13213>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach,



- H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, pp. 1–14, 2018.
- Alexandre Piché, Valentin Thomas, Cyril Ibrahim, Yoshua Bengio, and Chris Pal. Probabilistic planning with sequential monte carlo methods. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=ByetGn0cYX>.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.
- Siddharth Reddy, Anca D. Dragan, and Sergey Levine. {SQIL}: Imitation learning via reinforcement learning with sparse rewards. In *International Conference on Learning Representations (ICLR)*, 2020. URL <https://openreview.net/forum?id=S1xKd24twB>.
- Desik Rengarajan, Gargi Vaidya, Akshay Sarvesh, Dileep Kalathil, and Srinivas Shakkottai. Reinforcement learning with sparse rewards using guidance from offline demonstration. In *International Conference on Learning Representations (ICLR)*, 2022. URL <https://openreview.net/forum?id=YJ1WzgMVsmT>.
- Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. *Journal of machine learning research*, pp. 627–635, 2011.
- Dorsa Sadigh, Shankar Sastry, Sanjit A Seshia, and Anca D Dragan. Planning for autonomous cars that leverage effects on human actions. In *Robotics: Science and Systems XII, University of Michigan, Ann Arbor, Michigan, USA, June 18 - June 22, 2016*, 2016.
- Dorsa Sadigh, Anca D Dragan, Shankar Sastry, and Sanjit A Seshia. *Active preference-based learning of reward functions*. University of California, Berkeley, 2017.
- Stefan Schaal and Christopher G Atkeson. Learning control in robotics. In *2010 IEEE Robotics & Automation Magazine*, pp. 20–29. IEEE, 2010.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the 32th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, pp. 107–194, 2012.
- Lior Shani, Tom Zahavy, and Shie Mannor. Online apprenticeship learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 8240–8248, 2022.
- Wen Sun, Arun Venkatraman, Geoffrey J Gordon, Byron Boots, and J Andrew Bagnell. Deeply aggravated: Differentiable imitation learning for sequential prediction. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 3309–3318. PMLR, 2017.
- Wen Sun, J. Andrew Bagnell, and Byron Boots. Truncated horizon policy search: Deep combination of reinforcement and imitation. In *International Conference on Learning Representations (ICLR)*, 2018. URL <https://openreview.net/forum?id=ryUlhWCZ>.
- Wen Sun, Anirudh Vemula, Byron Boots, and Drew Bagnell. Provably efficient imitation learning from observation alone. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 6036–6045. PMLR, 2019.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- Umar Syed and Robert E Schapire. A game-theoretic approach to apprenticeship learning. In *Advances in neural information processing systems*, 2007. URL <https://proceedings.neurips.cc/paper/2007/file/ca3ec598002d2e7662e2ef4bdd58278b-Paper.pdf>.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE International Conference on Intelligent Robots and Systems (IROS)*, pp. 5026–5033. IEEE, 2012.
- Hado van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, and Joseph Modayil. Deep reinforcement learning and the deadly triad, 2018. URL <https://arxiv.org/abs/1812.02648>.
- Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *The 22nd international conference on artificial intelligence and statistics*, pp. 1195–1204. PMLR, 2019a. URL <https://arxiv.org/pdf/1810.07288.pdf>.
- Sharan Vaswani, Aaron Mishkin, Issam Laradji, Mark Schmidt, Gauthier Gidel, and Simon Lacoste-Julien. Painless stochastic gradient: Interpolation, line-search, and convergence rates. In *Advances in neural information processing systems*, 2019b. URL <https://proceedings.neurips.cc/paper/2019/file/2557911c1bf75c2b643afb4ecbfc8ec2-Paper.pdf>.
- Sharan Vaswani, Issam Laradji, Frederik Kunstner, Si Yi Meng, Mark Schmidt, and Simon Lacoste-Julien. Adaptive gradient methods converge faster with over-parameterization (but you should do a line-search), 2020. URL <https://arxiv.org/abs/2006.06835>.
- Andrew Warrington, Jonathan W Lavington, Adam Scibior, Mark Schmidt, and Frank Wood. Robust asymmetric learning in pomdps. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, pp. 11013–11023. PMLR, 2021. URL <https://proceedings.mlr.press/v139/warrington21a.html>.
- Yuege Xie, Xiaoxia Wu, and Rachel Ward. Linear convergence of adaptive stochastic gradient descent. In *Twenty Third International Conference on Artificial Intelligence and Statistics (AISTATS)*, Proceedings of Machine Learning Research, pp. 1475–1485. PMLR, 2020. URL <https://proceedings.mlr.press/v108/xie20a.html>.
- Zhaoming Xie, Glen Berseth, Patrick Clary, Jonathan Hurst, and Michiel van de Panne. Feedback control for cassie with deep reinforcement learning. In *2018 IEEE International Conference on Intelligent Robots and Systems (IROS)*, pp. 1241–1246. IEEE, 2018.
- Xinyan Yan, Byron Boots, and Ching-An Cheng. Explaining fast improvement in online policy optimization. In *Computing Research Repository*, 2020. URL <https://arxiv.org/abs/2007.02520>.
- Tom Zahavy, Alon Cohen, Haim Kaplan, and Yishay Mansour. Apprenticeship learning via frank-wolfe. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 6720–6728, 2020.
- Tom Zahavy, Brendan O’Donoghue, Guillaume Desjardins, and Satinder Singh. Reward is enough for convex mdps. In *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/pdf?id=ELndVeVA-TR>.
- Junyu Zhang, Alec Koppel, Amrit Singh Bedi, Csaba Szepesvari, and Mengdi Wang. Variational policy gradient method for reinforcement learning with general utilities. In *Advances in Neural Information Processing Systems*, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/30ee748d38e21392de740e2f9dc686b6-Paper.pdf>.
- Xinyi Zhang and Michiel van de Panne. Data-driven autocompletion for keyframe animation. In *Proceedings of the 11th Annual International Conference on Motion, Interaction, and Games*. Association for Computing Machinery, 2018. URL <https://doi.org/10.1145/3274247.3274502>.
- Sean X Zhou, Zhijie Tao, and Xiuli Chao. Optimal control of inventory systems with multiple types of remanufacturable products. *Manufacturing & Service Operations Management*, pp. 20–34, 2011.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (ICML)*, Proceedings of Machine Learning Research, pp. 928–936. PMLR, 2003. URL <https://www.aaai.org/Papers/ICML/2003/ICML03-120.pdf>.

---

## Supplementary material

---

### ORGANIZATION OF THE APPENDIX

- [A Definitions](#)
- [B Proof of Theorem 3.1](#)
- [C Proofs for Section 4](#)
- [D Additional Experimental Details](#)

### A DEFINITIONS

Our main assumptions are that each individual function  $l_i$  is differentiable, has a finite minimum  $l_i^*$ , and is  $L$ -smooth, meaning that for all  $v$  and  $w$ ,

$$l_i(v) \leq l_i(w) + \langle \nabla l_i(w), v - w \rangle + \frac{L}{2} \|v - w\|^2, \quad (\text{Individual Smoothness})$$

which also implies that  $f$  is  $L$ -smooth. A consequence of smoothness is the following bound on the norm of the stochastic gradients,

$$\|\nabla l_i(w)\|^2 \leq 2L(l_i(w) - l_i^*).$$

We also assume that each  $l_i$  is convex, meaning that for all  $v$  and  $w$ ,

$$l_i(v) \geq l_i(w) - \langle \nabla l_i(w), w - v \rangle, \quad (\text{Convexity})$$

Depending on the setting, we will also assume that  $f$  is  $\mu$  strongly-convex, meaning that for all  $v$  and  $w$ ,

$$f(v) \geq f(w) + \langle \nabla f(w), v - w \rangle + \frac{\mu}{2} \|v - w\|^2, \quad (\text{Strong Convexity})$$

### B PROOF OF THEOREM 3.1

**Theorem 3.1** (Follow-the-leader - Online Imitation Learning). *Under the following assumptions: (i) the policy class  $\Pi$  is sufficiently expressive so as to contain the expert policy, the (ii) optimization problem in Eq. (4) can be solved exactly, and (iii) the divergence  $D(\pi, \pi_e)$  has a unique minimizer and is bounded in the sufficient statistics of  $\pi$ , FTL (Eq. (4)) obtains the following regret guarantee in the OIL setting,*

$$R(T) \leq \frac{C}{1 - \gamma} \quad \text{where, } C := \max_{\tau, t} \left\{ \mathbb{E}_{s \sim p_t^\tau} [D(\pi(\cdot|s), \pi_e(\cdot|s))] \right\} < \max_s D(\pi(\cdot|s), \pi_e(\cdot|s)).$$

where  $p_t^\tau(s)$  is the probability of reaching state  $s$  at time-step  $\tau$  using policy  $\pi_t$ .

*Proof.* Recall that at every round FTL returns the following set of parameters:

$$\pi_T = \arg \min_{\pi} F_{T-1}(\pi) = \arg \min_{\pi} \sum_{t=1}^{T-1} l_t(\pi) = \arg \min_{\pi} \sum_{t=1}^{T-1} \mathbb{E}_{d_{\pi_t}^\gamma} [D(\pi, \pi_e)] \quad (9)$$

The loss at round  $t$  can be decomposed using the marginal state distribution at round  $t$ . Specifically, if  $p_{t-1}^\tau(s)$  is the probability of visiting state  $s$  at time-stamp  $\tau$  under the policy  $\pi_{t-1}$ , then,

$$l_t(\pi_t) = (1 - \gamma) \sum_{\tau=t}^{\infty} \gamma^\tau \mathbb{E}_{p_t^\tau} [D(\pi_t(\cdot|s), \pi_e(\cdot|s))] + (1 - \gamma) \sum_{\tau=0}^{t-1} \gamma^\tau \mathbb{E}_{p_t^\tau} [D(\pi_t(\cdot|s), \pi_e(\cdot|s))] \quad (10)$$

From Lemma B.1, we get that for all  $\tau < t$  we have  $p_t^\tau = p_{\pi_e}^\tau$ . This implies:

$$l_t(\pi_t) = (1 - \gamma) \sum_{\tau=t}^{\infty} \gamma^\tau \mathbb{E}_{\frac{p_t^\tau}{p_t^\tau}} [D(\pi_t(\cdot|s), \pi_e(\cdot|s))] + (1 - \gamma) \sum_{\tau=0}^{t-1} \gamma^\tau \mathbb{E}_{\frac{p_t^\tau}{p_{\pi_e}^\tau}} [D(\pi_t(\cdot|s), \pi_e(\cdot|s))] \quad (11)$$

In settings where we are able to exactly match the expert, by using FTL,  $\pi_t(\cdot|s) = \pi_e(\cdot|s)$  for all states where  $p_t^\tau(s) > 0$  for any  $i = 1, \dots, t-1$  and  $\tau < t$ . Using Lemma B.1, for  $\tau < i$ , since  $p_t^\tau = p_{\pi_e}^\tau$ , we can conclude that  $\pi_t(\cdot|s) = \pi_e(\cdot|s)$  for all states where  $p_{\pi_e}^\tau(s) > 0$  for  $\tau < t$ . This means that the second term  $\mathbb{E}_{p_{\pi_e}^\tau} [D(\pi_t(\cdot|s), \pi_e(\cdot|s))] = 0$ .

$$= (1 - \gamma) \sum_{\tau=t+1}^{\infty} \gamma^\tau \mathbb{E}_{\frac{p_t^\tau}{p_t^\tau}} [D(\pi(\cdot|s), \pi_e(\cdot|s))] \quad (12)$$

$$\leq \max \left\{ \mathbb{E}_{\frac{p_t^\tau}{p_t^\tau}} [D(\pi(\cdot|s), \pi_e(\cdot|s))] \right\} (1 - \gamma) \sum_{\tau=t+1}^{\infty} \gamma^\tau \quad (13)$$

$$\leq C \gamma^{t+1} \quad (14)$$

where  $C := \max \{ \mathbb{E}_{p_t^\tau} [D(\pi(\cdot|s), \pi_e(\cdot|s))] \}$ . We can now sum the left and right hand sides over  $T$ .

$$\sum_{t=1}^T l_t(\pi_t) - \sum_{t=1}^T l_t(\pi^*) \leq C \frac{1 - \gamma^{T+1}}{1 - \gamma} \leq \frac{C}{1 - \gamma} \quad (15)$$

□

**Lemma B.1** (Follow the Leader: Induced State-Distribution Under Interpolation). *Let the sequence of loss functions observed be defined as the following:*

$$l_t(\pi) = \mathbb{E}_{\frac{p_t^\tau}{p_t^\tau}} [D(\pi(\cdot|s), \pi_e(\cdot|s))], \quad l_t^\tau(\pi) = \int_s p_t^\tau(s) D(\pi(\cdot|s), \pi_e(\cdot|s)) ds \quad (16)$$

Assume that the policy class contains  $\pi_e$ . If at every round we minimize the following function with respect to  $\pi$ :

$$F_t(\pi) = \sum_{k=0}^{t-1} l_k(\pi) = (1 - \gamma) \sum_{k=0}^{t-1} \sum_{\tau=0}^{\infty} \gamma^\tau l_k^\tau(\pi) \quad (17)$$

then  $\forall \tau < t$  we have that  $p_t^\tau(s) = p_e^\tau(s)$ .

*Proof.* We now show inductively that at round  $t$ ,  $p_t(s) := p_{\pi_t}(s) = p_{\pi_e}(s) \quad \forall \tau < t$  under the assumptions stated above. The proof proceeds by induction, starting with our base case.

**Base-Case:**  $\tau = 0$

We can immediately note that the initial state distribution is independent of policy, and thus following the definition of our generative model  $p_1^0(s_0) = p_e^0(s_0)$ .

**Base-Case:**  $\tau = 1$

Next we can consider the results of a single round of the algorithm, and its effect on  $p_1^1(s)$ . By following FTL under interpolation we have:

$$\pi_1 = \arg \min_{\pi} L_1(\pi) = \arg \min_{\pi} l_0(\pi) \quad (\text{Definition of } L_1.)$$

$$l_0(\pi_1) = 0 \implies l_0^k(\pi_1) = 0 \quad \forall k. \quad (\text{Interpolation assumption.})$$

$$\implies D[\pi_1(\cdot|s), \pi_e(\cdot|s)] = 0 \quad \forall s \text{ s.t. } p_0^0(s) := p_e^0(s) > 0 \quad (\text{Definition of } l_0^1.)$$

$$\implies \pi_1(a|s) = \pi_e(a|s) \quad \forall s \text{ s.t. } p_0^0(s) > 0, \quad (\text{Definition of a divergence and assumption stated below.})$$

The final line holds provided the divergence  $D$  is strongly convex with respect to the policy  $\pi$ . Now we have that both  $p_1^0(s) = p_e^0(s)$  by  $\tau=0$ , and  $\pi_1(a|s) = \pi_e(a|s)$ , we can get direct equality with respect to state at

the next time-step:

$$p_e^1(s') = \int_s \int_a p(s'|s, a) \pi_e(a|s) p_e^0(s) ds da \quad (18)$$

$$= \int_s \int_a p(s'|s, a) \pi_1(a|s) p_0^0(s) ds da \quad (19)$$

$$= p_1^1(s') \quad (20)$$

### Inductive Step:

Assume that for some arbitrary round  $t$  we have the following:  $p_e^\tau(s) = p_{t-1}^\tau(s) \quad \forall \tau < t$ , and we want to prove that at round  $t+1$ ,  $p_e^\tau(s) = p_t^\tau(s) \quad \forall \tau < t+1$ ,

(1) Following the same argument as before  $\pi_t(a|s) = \pi_e(a|s) \quad \forall a, \quad \forall s$  s.t.  $p_e^{t-1}(s) > 0$ :

(2) We can show that  $p_{t-1}^{t-1}(s) = p_t^{t-1}(s) \quad \forall s_{t-1}$  s.t.  $p_e^{t-1}(s) > 0$  by way of contradiction. Consider the case where the above does not hold, but the policy interpolates all previous functions observed. By definition of  $\pi_t$  and  $\pi_{t-1}$  for all  $\tau < t$  if  $p_e^\tau(s) > 0$  then  $\pi_t(\cdot|s) = \pi_e(\cdot|s)$  and  $\pi_{t-1}(\cdot|s) = \pi_e(\cdot|s)$ .

If  $p_{t-1}^{t-1}(s) \neq p_t^{t-1}(s) \quad \forall s_{t-1}$  s.t.  $p_e^{t-1}(s) > 0$ , then let us consider the first time step  $\tau^*$  such that  $p_{t-1}^{\tau^*}(s) \neq p_t^{\tau^*}(s) \quad \forall s$  s.t.  $p_{t-1}^{\tau^*}(s) = p_e^{\tau^*}(s) > 0$ . Then by definition of  $\tau^*$  and the inductive hypothesis,  $p_{t-1}^{\tau^*-1}(s) = p_t^{\tau^*-1}(s) \quad \forall s$  s.t.  $p_{t-1}^{\tau^*-1}(s) = p_e^{\tau^*-1}(s) > 0$ . Now by the definition of the generative model this implies that  $\exists s$  s.t.  $\pi_t(\cdot|s) \neq \pi_{t-1}(\cdot|s)$  and  $p_e^{\tau^*-1}(s) > 0$ . This is a contradiction, as we showed above. Hence  $p_{t-1}^{t-1}(s) = p_t^{t-1}(s) \quad \forall s_{t-1}$  s.t.  $p_e^{t-1}(s) > 0$ .

We can again put all this together to get the desired result:

$$\begin{aligned} p_t^t(s') &= \int_s \int_a p(s'|s, a) \pi_t(a|s) p_t^{t-1}(s) ds da && \text{(Definition of Marginal.)} \\ &= \int_s \int_a p(s'|s, a) \pi_t(a|s) p_{t-1}^{t-1}(s) ds da && \text{(Sub-proof (2).)} \\ &= \int_s \int_a p(s'|s, a) \pi_e(a|s) p_e^{t-1}(s) ds da && \text{(Inductive hypothesis.)} \\ &= \int_s \int_a p(s'|s, a) \pi_e(a|s) p_e^t(s) ds da && \text{(FTL)} \\ &= p_e^t(s'). && \text{(Definition of Marginal.)} \end{aligned}$$

Therefore by induction we have that under interpolation, FTL ensures that at round  $t$   $p_t^\tau(s) = p_e^\tau(s) \quad \forall \tau \leq t-1$ .  $\square$

## C PROOFS FOR SECTION 4

### C.1 PROOF OF PROPOSITION 4.1

**Proposition 4.1** (Reformulation). *Defining  $\eta_t := 1/(\sum_{i=1}^t \sigma_i)$ , the update in Eq. (5) can be reformulated (proof in Appendix C) as:*

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} \left[ \sum_{i=1}^t l_i(w) - \langle w, \sum_{i=1}^{t-1} \nabla l_i(w_t) \rangle + \frac{1}{2\eta_t} \|w - w_t\|^2 \right]. \quad (6)$$

*Proof.* Since  $\eta_t := \frac{1}{\sum_{s=1}^t \sigma_s}$ , by definition of  $w_{t+1}$  in Eq. (5),

$$\sum_{s=1}^t \nabla l_s(w_{t+1}) + \frac{w_{t+1}}{\eta_t} = \sum_{s=1}^t \sigma_s w_s$$



Similarly, by definition of  $w_t$ ,

$$\sum_{s=1}^{t-1} \nabla l_s(w_t) + \frac{w_t}{\eta_{t-1}} = \sum_{s=1}^{t-1} \sigma_s w_s$$

From the above relations,

$$\begin{aligned} \sum_{s=1}^t \nabla l_s(w_{t+1}) + \frac{w_{t+1}}{\eta_t} &= \sum_{s=1}^{t-1} \nabla l_s(w_t) + \sigma_t w_t + \frac{w_t}{\eta_{t-1}} \\ \implies w_{t+1} + \eta_t \sum_{s=1}^t \nabla l_s(w_{t+1}) &= w_t + \eta_t \sum_{s=1}^{t-1} \nabla l_s(w_t) \end{aligned}$$

Therefore, at iteration  $t$ , we need to obtain  $w_{t+1}$ , we need to solve the following equation w.r.t  $w$ ,

$$w + \eta_t \sum_{s=1}^t \nabla l_s(w) = w_t + \eta_t \sum_{s=1}^{t-1} \nabla l_s(w_t) \quad (21)$$

Similar to FTL, this update requires storing the previous functions  $f_s$  from  $s = 1$  to  $t$ , but does not require storing the previous models like in Eq. (5). Minimizing the following loss is equivalent to ensuring Eq. (21).

$$w_{t+1} = \arg \min_w \left[ \sum_{s=1}^t l_s(w) - \langle w, \sum_{s=1}^{t-1} \nabla l_s(w_t) \rangle + \frac{1}{2\eta_t} \|w - w_t\|^2 \right] \quad (22)$$

□

## C.2 DERIVATION OF ALTERNATIVE FTRL REFORMULATION IN EQUATION 7

In a more direct fashion than the proof above, we can see that

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} \left[ \sum_{i=1}^t l_i(w) + \frac{1}{2\eta_t} \|w\|^2 - \sum_{i=1}^{t-1} w^\top w_i \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] \right], \quad (23)$$

In fact represents the same objective as

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} \left[ \sum_{i=1}^t l_i(w) + \sum_{i=1}^t \frac{\sigma_i}{2} \|w - w_i\|^2 \right] \quad (24)$$

We can see this if we differentiate both equations

$$\begin{aligned} \nabla F_{\text{alt}}(w) &= \nabla \left[ \sum_{i=1}^t l_i(w) + \frac{1}{2\eta_t} \|w\|^2 - \sum_{i=1}^{t-1} w^\top w_i \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] \right] \\ &= \nabla \left[ \sum_{i=1}^t l_i(w) + \frac{w^\top w}{2} \sum_{i=1}^{t-1} \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] - \sum_{i=1}^{t-1} w^\top w_i \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] \right] \\ &= \nabla \left[ \sum_{i=1}^t l_i(w) + \frac{w^\top w}{2} \sum_{i=1}^{t-1} \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] - \sum_{i=1}^{t-1} w^\top w_i \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] + \sum_{i=1}^{t-1} \frac{w_i^\top w_i}{2} \left[ \frac{1}{\eta_t} - \frac{1}{\eta_{i-1}} \right] \right] \\ &= \nabla \left[ \sum_{i=1}^t l_i(w) + \sum_{i=1}^{t-1} \frac{\sigma_i}{2} \|w - w_i\|^2 \right] \\ &= \nabla F(w) \end{aligned} \quad (25)$$

Note that we define,

$$\sigma_t = \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}}, \quad \text{and} \quad \frac{1}{\eta_t} = \sqrt{t} \quad (26)$$

to ensure that the magnitude of regularization used in both reformulations is of order  $\sqrt{t}$ .

## C.3 MAIN REGRET LEMMA

We define

$$F_t(w) := \sum_{i=0}^{t-1} l_i(w) + \psi_t(w), \quad (27)$$

where,  $\psi_t(w)$  is a strongly-convex proximal regularizer that satisfies the following property:

$$w_t = \arg \min [\psi_{t+1}(w) - \psi_t(w)]. \quad (28)$$

Eq. (5) uses  $\psi_t(w) = \sum_{i=1}^{t-1} \frac{\sigma_i}{2} \|w - w_i\|^2$ . Since  $\psi_{t+1}(w) - \psi_t(w) = \frac{\sigma_t}{2} \|w - w_t\|^2$ , which is minimized at  $w_t$  and hence the regularizer in Eq. (5) satisfies the desired property, and is a valid strongly-convex proximal regularizer. We will now prove the following lemma for a general strongly-convex proximal regularizer  $\psi_t$ . In this case, the FTRL update in Eq. (5) can be generalized to:

$$w_t = \arg \min F_t(w) = \arg \min \sum_{i=0}^{t-1} l_i(w) + \psi_t(w) \implies \sum_{i=1}^{t-1} \nabla l_i(w_t) + \nabla \psi_t(w_t) = 0. \quad (29)$$

**Lemma C.1.** *Assuming that the functions  $F_t$  are  $\lambda_t$ -strongly convex, then the regret for the FTRL update in Eq. (29) can be bounded as:*

$$R(T) \leq [F_1(w_1)] + \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] - \sum_{t=1}^T [\psi_t(w_t) - \psi_{t+1}(w_t)] + \psi_{T+1}(w^*)$$

*Proof.*

$$\begin{aligned} F_{t+1}(w_t) - F_{t+1}(w_{t+1}) &\leq \frac{1}{2\lambda_{t+1}} \|\nabla F_{t+1}(w_t)\|^2 \\ &\quad \text{(Since } F_{t+1} \text{ is } \lambda_{t+1}\text{-strongly convex and } w_{t+1} \text{ is the minimizer of } F_{t+1}\text{.)} \\ &= \frac{1}{2\lambda_{t+1}} \left\| \sum_{i=1}^t \nabla l_i(w_t) + \nabla \psi_{t+1}(w_t) \right\|^2 \quad \text{(By definition of } F_{t+1}\text{)} \\ &= \frac{1}{2\lambda_{t+1}} \left\| \sum_{i=1}^{t-1} \nabla l_i(w_t) + \nabla \psi_t(w_t) + \nabla \psi_{t+1}(w_t) - \nabla \psi_t(w_t) + \nabla l_t(w_t) \right\|^2 \\ &= \frac{1}{2\lambda_{t+1}} \|\nabla F_t(w_t) + \nabla \psi_{t+1}(w_t) - \nabla \psi_t(w_t) + \nabla l_t(w_t)\|^2 \quad \text{(By definition of } F_t\text{)} \\ &= \frac{1}{2\lambda_{t+1}} \|\nabla \psi_{t+1}(w_t) - \nabla \psi_t(w_t) + \nabla l_t(w_t)\|^2 \quad \text{(Since } w_t \text{ is the minimizer of } F_t\text{)} \\ &= \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \quad \text{(Since } w_t \text{ is the minimizer of } \psi_{t+1}(w) - \psi_t(w)\text{)} \end{aligned}$$

$$\begin{aligned} F_{t+1}(w_t) - F_{t+1}(w_{t+1}) &= [F_{t+1}(w_t) - F_t(w_t)] + [F_t(w_t) - F_{t+1}(w_{t+1})] \\ &= [l_t(w_t) + \psi_{t+1}(w_t) - \psi_t(w_t)] + [F_t(w_t) - F_{t+1}(w_{t+1})] \end{aligned}$$

Summing from  $t = 1$  to  $T$ , and using the above relation,

$$\begin{aligned} \sum_{t=1}^T [l_t(w_t) + \psi_{t+1}(w_t) - \psi_t(w_t)] + \sum_{t=1}^T [F_t(w_t) - F_{t+1}(w_{t+1})] &\leq \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] \\ \implies \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] + [F_1(w_1) - F_{T+1}(w_{T+1})] &\leq \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] - \sum_{t=1}^T [\psi_t(w_t) - \psi_{t+1}(w_t)] - \sum_{t=1}^T l_t(w^*) \\ R(T) \leq [F_{T+1}(w_{T+1}) - F_{T+1}(w^*) - F_1(w_1)] &+ \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] - \sum_{t=1}^T [\psi_t(w_t) - \psi_{t+1}(w_t)] + \psi_{T+1}(w^*) \end{aligned}$$

Since  $w_{T+1}$  is the minimizer of  $F_{T+1}$ ,

$$R(T) \leq [F_1(w_1)] + \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] - \sum_{t=1}^T [\psi_t(w_t) - \psi_{t+1}(w_t)] + \psi_{T+1}(w^*)$$

□

The above expression is true for both FTL, and (adaptive) FTRL, and only uses the definitions of the proximal regularizer and the strong-convexity property for  $F_t$ . We specialize this result for  $\psi_t(w) = \sum_{i=1}^{t-1} \frac{\sigma_i}{2} \|w - w_i\|^2$  used in Eq. (5).

**Lemma C.2.** *Assuming that each  $l_i$  is  $\mu_i$  strongly-convex for  $\mu_i \geq 0$ , the regret for the FTRL update in Eq. (5) can be bounded as:*

$$R(T) \leq \sum_{t=1}^T \left[ \frac{1}{2 \sum_{i=1}^t [\sigma_i + \mu_i]} \|\nabla l_t(w_t)\|^2 \right] + \frac{D^2}{2} \sum_{t=1}^T \sigma_t$$

where  $D$  is the diameter of  $\mathcal{W}$ .

*Proof.* With this choice of  $\psi_t$ , we note that  $F_1(x_1) = \psi_1(x) = 0$  for all  $w$ . Using Lemma C.1,

$$\begin{aligned} R(T) &\leq \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] + \sum_{t=1}^T \left[ \frac{\sigma_t}{2} \|w_t - w^*\|^2 \right] + \psi_{T+1}(w^*) \\ &= \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] + \sum_{t=1}^T \frac{\sigma_t}{2} \|w^* - w_t\|^2. \end{aligned}$$

Since the iterates are bounded on  $\mathcal{W}$ ,  $\|w - w^*\| \leq D$  for all  $w$ ,

$$R(T) \leq \sum_{t=1}^T \left[ \frac{1}{2\lambda_{t+1}} \|\nabla l_t(w_t)\|^2 \right] + \frac{D^2}{2} \sum_{t=1}^T \sigma_t$$

Since each  $l_i$  is  $\mu_i$  strongly-convex,  $F_t$  is  $\sum_{i=1}^{t-1} \mu_i$  strongly-convex, and hence  $\lambda_{t+1} = \sum_{i=1}^t \mu_i$ .

□

#### C.4 PROOF OF THEOREM 4.3

**Theorem 4.3** (FTRL - Smooth + Convex). *Assuming each  $l_t$  is (i)  $L$ -smooth, (ii) convex, FTRL (Eq. (6)) for  $\eta_t = \min \left\{ \left( \sum_{i=1}^t \epsilon_i^2 \right)^{-1/2}, \frac{1}{2L} \right\}$  for all  $t$ , achieves the following regret*

$$R(T) \leq \sum_{t=1}^T \left[ \frac{\eta_t \|\nabla l_t(w_t)\|^2}{2} \right] + \frac{D^2}{2\eta_T} \leq 2D^2L + (D^2 + 2L) \sqrt{\sum_{t=1}^T \epsilon_t^2}.$$

*Proof.* In this case,  $\mu_i = 0$  i.e.  $l_i$  is convex only convex without strong-convexity. Using Lemma C.2 with  $\mu_i = 0$  and defining  $\eta_t := \frac{1}{\sum_{i=1}^t \sigma_i}$ ,

$$R(T) \leq \sum_{t=1}^T \left[ \frac{\eta_t}{2} \|\nabla l_t(w_t)\|^2 \right] + \frac{D^2}{2\eta_T} \quad (30)$$

Recall that we use a constant step-size implying that  $\eta_1 = \eta_2 = \eta_T = \eta = \min\{\frac{1}{\mathcal{E}}, \frac{1}{2L}\}$ , where  $\mathcal{E}^2 := \sum_{t=1}^T \epsilon_t^2$ . With this choice,

$$\begin{aligned} R(T) &\leq \eta L \sum_{t=1}^T [l_t(w_t) - l_t(w_t^*)] + \frac{D^2}{2\eta} && \text{(By smoothness, and since } w_t^* \text{ is a minimizer of } l_t.) \\ &= \eta L \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] + \eta L \sum_{t=1}^T [l_t(w^*) - l_t(w_t^*)] + \frac{D^2}{2\eta} \\ &= \eta L R(T) + \eta L + \sum_{t=1}^T \epsilon_t^2 + \frac{D^2}{2\eta} \end{aligned}$$

Since  $\eta < \frac{1}{L}$ ,

$$R(T) \leq \frac{\eta L}{1 - \eta L} \sum_{t=1}^T \epsilon_t^2 + \frac{D^2}{2\eta(1 - \eta L)}$$

Since  $\eta \leq \frac{1}{2L}$ ,  $\frac{1}{1 - \eta L} \leq 2$ ,

$$\implies R(T) \leq 2\eta L \sum_{t=1}^T \epsilon_t^2 + \frac{D^2}{\eta}$$

Since  $\eta = \min\{\frac{1}{\mathcal{E}}, \frac{1}{2L}\}$ ,  $\frac{1}{\eta} = \max\{\mathcal{E}, 2L\}$ ,

$$\begin{aligned} R(T) &\leq 2\eta L \sum_{t=1}^T \epsilon_t^2 + D^2 \max\{\mathcal{E}, 2L\} \leq \frac{2L}{\mathcal{E}} \sum_{t=1}^T \epsilon_t^2 + D^2(\mathcal{E} + 2L) \leq 2L\mathcal{E} + D^2\mathcal{E} + 2D^2L \\ \implies R(T) &\leq 2D^2L + (D^2 + 2L) \sqrt{\sum_{t=1}^T \epsilon_t^2} \end{aligned}$$

□

#### C.5 PROOF OF THEOREM 4.4

**Theorem 4.4** (AdaFTRL - Smooth + Convex). *Assuming each  $l_t$  is (i)  $L$ -smooth, (ii) convex, FTRL (Eq. (6)) for  $\eta_t = \alpha/\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}$ , achieves the following regret*

$$R(T) \leq \sum_{t=1}^T \left[ \frac{\alpha \|\nabla l_t(w_t)\|^2}{2 \sum_{i=1}^t \|\nabla l_i(w_i)\|^2} \right] + \frac{D^2}{2\alpha} \sqrt{\sum_{t=1}^T \|\nabla l_t(w_t)\|^2} \leq 2L \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right)^2 + \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T \epsilon_t^2}.$$

*Proof.* We follow the same proof as Theorem 4.3 until Eq. (30). Since  $\eta_t = \frac{\alpha}{\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}}$ ,

$$R(T) \leq \frac{\alpha}{2} \sum_{t=1}^T \left[ \frac{\|\nabla l_t(w_t)\|^2}{\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}} \right] + \frac{D^2}{2\alpha} \sqrt{\sum_{i=1}^T \|\nabla l_i(w_i)\|^2}$$

Bounding  $\sum_{t=1}^T \left[ \frac{\|\nabla l_t(w_t)\|^2}{\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}} \right] \leq \sqrt{\sum_{t=1}^T \|\nabla l_t(w_t)\|^2}$  using the AdaGrad inequality in (Duchi et al., 2011; Levy et al., 2018),

$$R(T) \leq \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T \|\nabla l_t(w_t)\|^2}$$

By smoothness, and since  $w_t^*$  is the minimizer of  $l_t$ .

$$\begin{aligned} &\leq \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T [l_t(w_t) - l_t(w_t^*)]}, \\ &= \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T [l_t(w_t) - l_t(w^*) + l_t(w^*) - l_t(w_t^*)]} \end{aligned}$$

Recall that  $\epsilon_t^2 := l_t(w^*) - l_t(w_t^*)$ , and using the definition of  $R(T)$ ,

$$\implies \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] \leq \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T [l_t(w_t) - l_t(w^*)] + \sum_{t=1}^T \epsilon_t^2}$$

Squaring both sides,

$$\left( \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] \right)^2 \leq 2L \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right)^2 \left( \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] + \sum_{t=1}^T \epsilon_t^2 \right)$$

Using Lemma C.3,

$$R(T) = \sum_{t=1}^T [l_t(w_t) - l_t(w^*)] \leq 2L \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right)^2 + \sqrt{2L} \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T \epsilon_t^2}$$

□

**Lemma C.3.** *If  $x^2 \leq a(x + b)$  for  $a \geq 0$  and  $b \geq 0$ ,*

$$x \leq \frac{1}{2} \sqrt{a^2 + 4ab} + a \leq a + \sqrt{ab}.$$

*Proof.* The starting point is the quadratic inequality  $x^2 - ax - ab \leq 0$ . Letting  $r_1 \leq r_2$  be the roots of the quadratic, the inequality holds if  $x \in [r_1, r_2]$ . The upper bound is then given by using  $\sqrt{a + b} \leq \sqrt{a} + \sqrt{b}$

$$r_2 = \frac{a + \sqrt{a^2 + 4ab}}{2} \leq \frac{a + \sqrt{a^2} + \sqrt{4ab}}{2} = a + \sqrt{ab}.$$

□

## C.6 FTRL IN THE NON-SMOOTH, BUT LIPSCHITZ SETTING

In the absence of smoothness, we will make the standard assumption that each  $l_i$  is  $G$ -Lipschitz, meaning that for all  $w$ ,  $\|\nabla l_i(w)\| \leq G$ .

**Theorem C.4** (FTRL - Lipschitz + Convex). *Assuming each  $l_t$  is (i)  $G$ -Lipschitz, (ii) convex, FTRL with  $\eta_t = \frac{\alpha}{\sqrt{t}}$  achieves the following regret,*

$$R(T) \leq \frac{\sqrt{T}}{2} \left[ G^2 \alpha + \frac{D^2}{\alpha} \right]$$

where  $D$  is the diameter of  $\mathcal{W}$ .

*Proof.* Using Lemma C.2 when  $\mu_i = 0$ , defining  $\eta_t := \frac{1}{\sum_{i=1}^t \sigma_t}$  and bounding  $\|\nabla l_t(w_t)\|^2 \leq G^2$ ,

$$R(T) \leq \frac{G^2}{2} \sum_{t=1}^T [\eta_t] + \frac{D^2}{2\eta_T}$$



For  $\eta_t = \frac{\alpha}{\sqrt{t}}$ ,

$$R(T) \leq \frac{G^2 \alpha \sqrt{T}}{2} + \frac{D^2 \sqrt{T}}{2\alpha}$$

□

**Theorem C.5** (AdaFTRL - Lipschitz + Convex). *Assuming each  $l_t$  is (i)  $G$ -Lipschitz, (ii) convex, AdaFTRL with  $\eta_t = \frac{\alpha}{\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}}$  achieves the following regret,*

$$R(T) \leq \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) G \sqrt{T}$$

where  $D$  is the diameter of  $\mathcal{W}$ .

*Proof.* We follow the same proof as Theorem 4.4. Since  $\eta_t = \frac{\alpha}{\sqrt{\sum_{i=1}^t \|\nabla l_i(w_i)\|^2}}$ ,

$$R(T) \leq \left( \frac{\alpha}{2} + \frac{D^2}{2\alpha} \right) \sqrt{\sum_{t=1}^T \|\nabla l_t(w_t)\|^2}$$

Bounding  $\|\nabla l_t(w_t)\|^2 \leq G^2$  completes the proof. □

#### C.7 FTL IN THE SMOOTH, STRONGLY-CONVEX SETTING

**Theorem C.6** (FTL - Smooth + Convex). *Assuming that (i) each  $l_i$  is  $\mu$  strongly-convex for  $\mu > 0$ , (ii) smooth, the regret for the FTL update in Eq. (4) can be bounded as:*

$$R(T) \leq \frac{DL}{\mu} (1 + \log(T))$$

*Proof.* FTL can be considered as a special case of the general FTRL update in Eq. (29) with  $\psi_t(w) = 0$  for all  $t$  and  $w$ , meaning that  $\sigma_t = 0$ . Using Lemma C.2 in this case,

$$R(T) \leq \sum_{t=1}^T \left[ \frac{1}{2 \sum_{i=1}^t [\mu]} \|\nabla l_t(w_t)\|^2 \right]$$

Using smoothness and since  $\mathcal{W}$  has diameter  $D$ ,

$$\leq DL \sum_{t=1}^T \frac{1}{2\mu t} \leq \frac{DL}{\mu} (1 + \log(T)).$$

□

#### C.8 FTL IN THE STRONGLY-CONVEX, NON-SMOOTH, BUT LIPSCHITZ SETTING

**Theorem C.7** (FTL - Lipschitz + Convex). *Assuming that (i) each  $l_i$  is  $\mu$  strongly-convex for  $\mu > 0$ , (ii)  $G$ -Lipschitz, the regret for the FTL update in Eq. (4) can be bounded as:*

$$R(T) \leq \frac{G^2}{2\mu} (1 + \log(T))$$

*Proof.* Following the same proof as Theorem C.6,

$$R(T) \leq \sum_{t=1}^T \left[ \frac{1}{2\mu t} \|\nabla l_t(w_t)\|^2 \right]$$

Since  $\|\nabla l_t(w_t)\|^2 \leq G^2$ ,

$$R(T) \leq \frac{G^2}{2\mu}(1 + \log(T)).$$

□

## D ADDITIONAL EXPERIMENTAL DETAILS

### D.1 GRID-WORLD EXAMPLE IN SECTION 1

In this experiment, we compare three algorithms: online gradient decent, follow-the-leader, and follow-the-regularized-leader. In this setting we create a  $7 \times 7$  grid where the agent is able to move in one of five ways at each time-step: up, down, left, right, or not at all. The expert actions switch in each round with up and right in odd rounds and down and left on even rounds Fig. 1. The experiment is run for a total of 100 rounds, each of which has a horizon of 5 environment interactions. At the beginning of each round, the agent starts at a position sampled uniformly at random. We verify that FTL incurs substantially larger regret over 1000 rounds. In Fig. 1, we only include the first 100 rounds to reduce the computational burden.

To solve the subproblems required for the FTL and the FTRL updates, we use an Armijo, backtracking line-search starting with a fixed step-size (Armijo, 1966). For OGD and FTRL, we use a grid-search over the  $\eta \in [10^{-5}, \dots, 10^5]$ . For a full list of parameters and code, see the *Gridworld-Example* folder inside the accompanying code-base. We also include a combined figure below illustrating the relationship between the reward, which is defined as zero if the agent did not chose the expert action, and one otherwise.

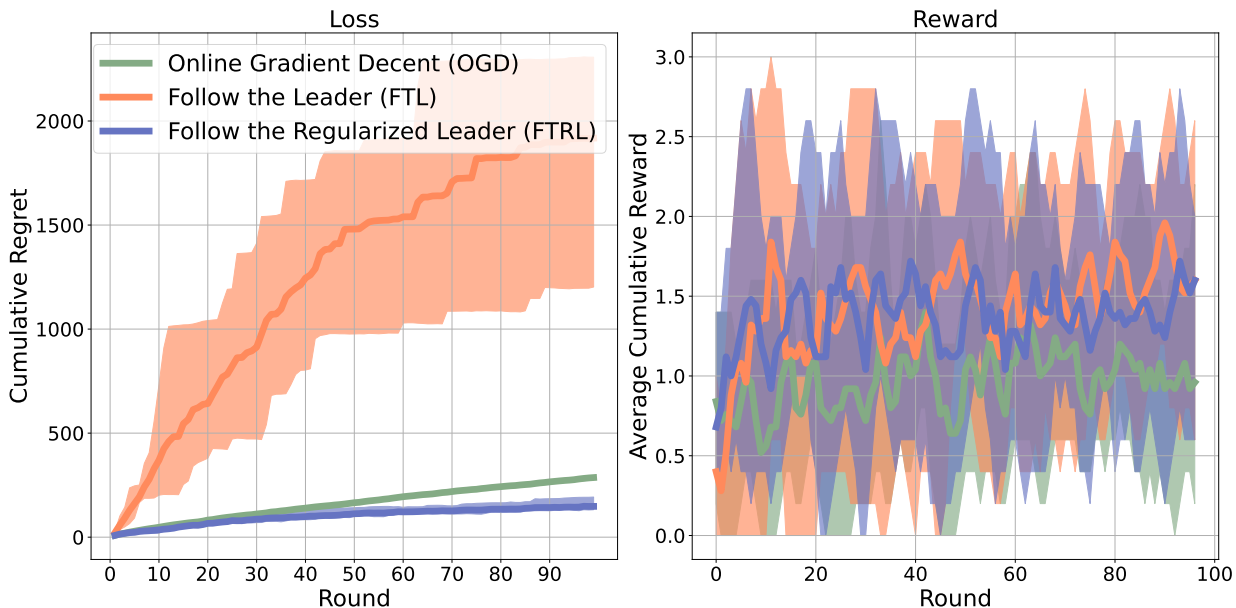


Figure 5: Displays the cumulative regret, and the cumulative reward over 100 rounds.

### D.2 TOY EXPERIMENTS

In these experiments, we consider a series of online regression problems which constitute a mixture of continuous state and discrete action problems (with a logistic loss), and continuous state and continuous action problems (with L1, L2 losses). For setting the hyper-parameters for these problems, we did a grid-search for each problem, loss class, and algorithm, for the initial ten rounds. We then used the best hyper-parameter configuration (in terms of the cumulative regret), and evaluated it for 250 rounds. Similar to Appendix D.1, we use a back-tracking line-search. We analyze the performance of each algorithm in two settings for each loss – *simple* and *adversarial*.

In the simple setting, we construct a random weight matrix  $W^* \in \mathbb{R}^{3 \times 10}$  and a random feature matrix  $X$  where each state is represented by a 10-dimensional vector. The action space consists of 3 actions generated as  $W^*X$ . For the discrete action case, the logits for each action are proportional to  $W^*X$ . In the adversarial setting, we use the same generative process as the simple setting, but switch between using  $W^{*\top}$  and  $-W^{*\top}$  in alternate rounds. For both settings, we sample one environment interaction/round. For a full list of parameters and code, see the *Toy-Experiments* folder inside the accompanying code-base.

The figures in Appendix D.2 display similar trends – (i) in the simple (non-adversarial) setting, across all losses, FTRL, AFTRL (AdaFTRL), FTL perform significantly better than AdaGrad, OGD (ii) in the adversarial setting, FTL has poor empirical performance across losses.

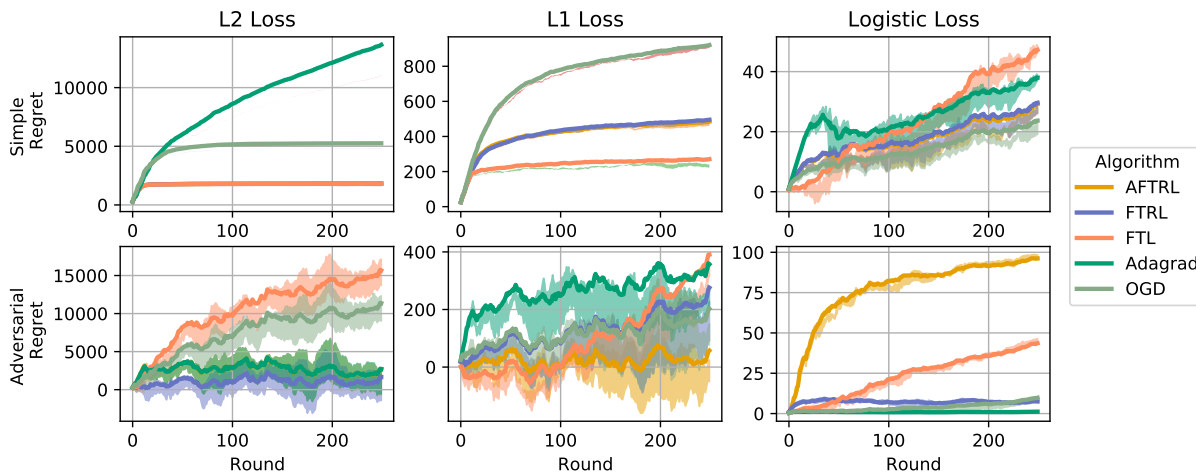


Figure 6: Displays a set of toy experiments in which different online learning algorithms are applied to various online regression problems. In this case we display both the L1 and L2 loss as well as the logistic loss. Our code also includes results for the Huber loss which were omitted.

### D.3 MUJOCO

In this section, we discuss the details for the Mujoco (Todorov et al., 2012) continuous control experiments. For a full list of parameters and code, see the *Mujoco-Experiments* folder inside the accompanying code-base.

#### D.3.1 BEHAVIORAL CLONING (INTERACTION UNDER EXPERT)

In the experiments presented in the main paper, we interact with the environment using only the current agent policy. Most existing algorithms (Ross et al., 2011) use a linear combination of the expert policy and the agents learned policy. We include an ablation in which we only use the expert policy to interact with the environment, and compare the performance of the algorithms for the linear and neural network settings in Fig. 7 and Fig. 8 respectively.

#### D.3.2 DECREASED SAMPLE-SIZES FOR LINEAR MODELS

In this section, we demonstrate the effect of using a reduced number of environment interactions per round. We use 100 environment interactions/round and show the results for the linear policy parameterization (Fig. 9, Fig. 10).

### D.4 ATARI

As was described in the main paper, in this setting both the expert and learned policy parameterize a categorical distribution which takes an input an 256 by 256 image of the atari screen. Following the same data augmentation as Schulman et al. (2015), we additionally convert the image to a 84 by 84 grey scale image and stack the previous four states observed (also called frame stacking). In this setting the expert is

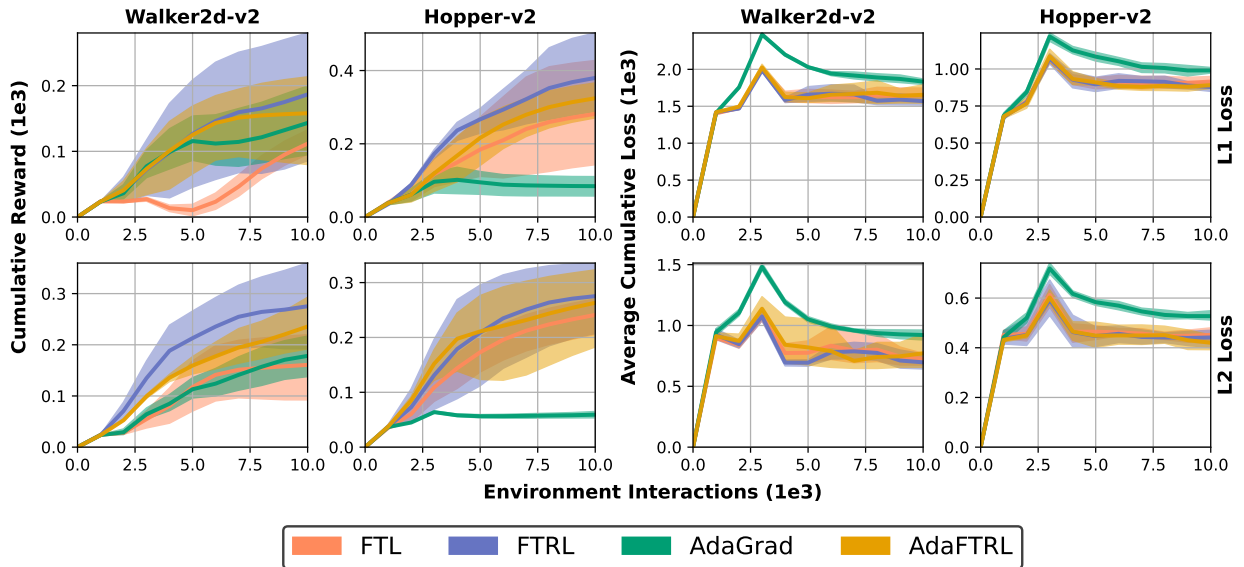


Figure 7: Comparing the algorithms with a linear policy parameterization when using the expert policy to interact with the environment. In this setting, we use 1000 environment interactions per round. We observe that FTRL, AdaFTRL are still the best performing methods, but FTL has slightly worse performance, with AdaGrad performs slightly better.

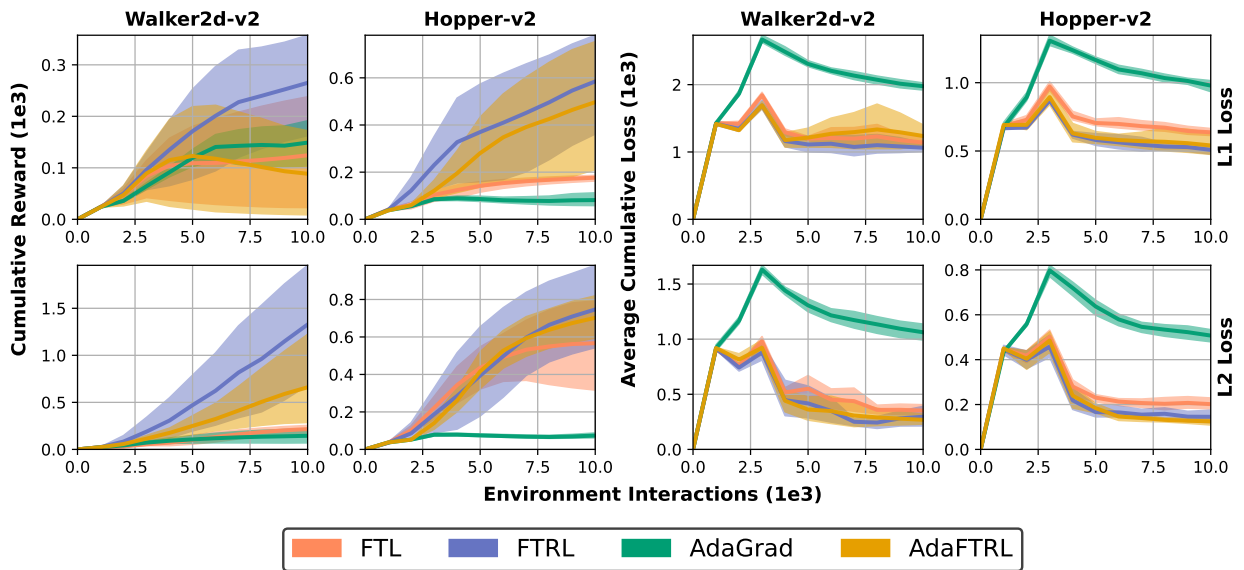


Figure 8: Comparing the algorithms with a neural network policy parameterization when using the expert policy to interact with the environment. In this setting, we use 1000 environment interactions per round. We observe that FTRL, AdaFTRL are still the best performing methods, but both FTL and AdaGrad have significantly worse performance.

learned using the PPO algorithm, and we again use the architecture described in [Schulman et al. \(2015\)](#). To reduce the computational burden, we only use 3 independent runs. For a full list of parameters and code, see the *Atari-Experiments* folder inside the accompanying code-base.

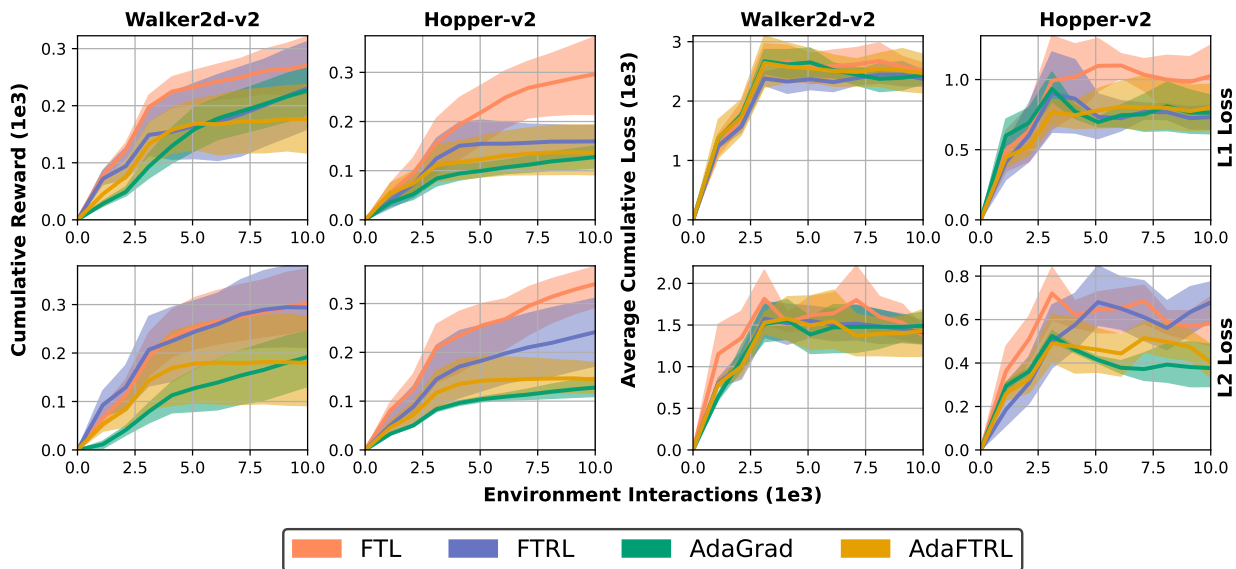


Figure 9: Comparing the algorithms with a linear policy parameterization when using the agent policy to interact with the environment. In this setting, we use 100 environment interactions per round. We observe that FTL, FTRL, AdaFTRL are the best performing methods, while AdaGrad performs better compared to Fig. 2.

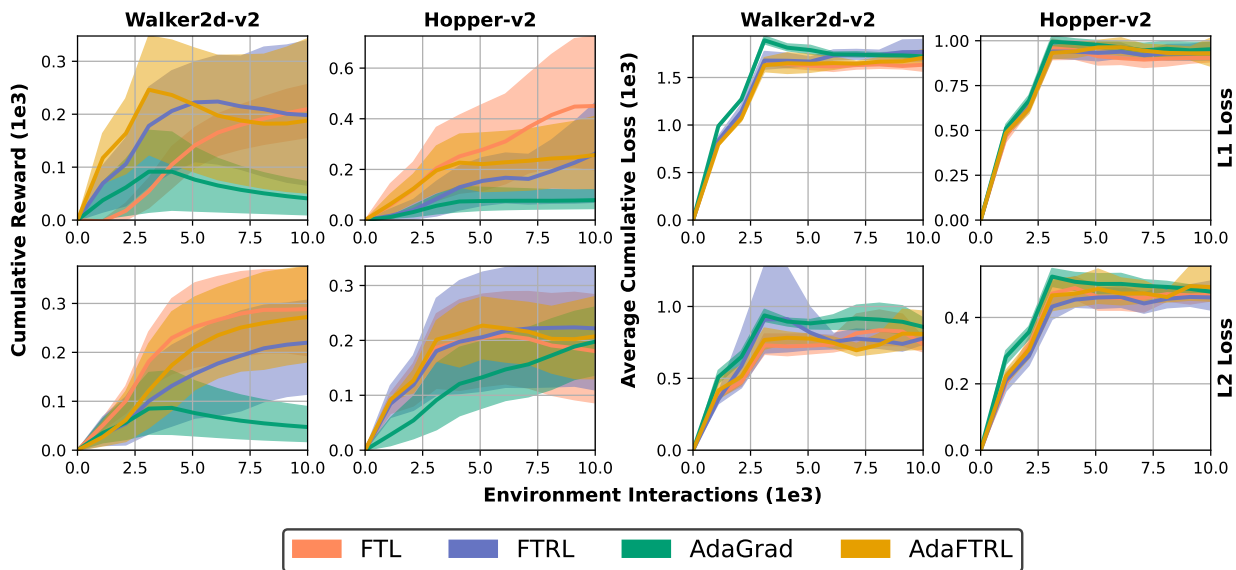


Figure 10: Comparing the algorithms with a linear policy parameterization when using the expert policy to interact with the environment. In this setting, we use 100 environment interactions per round. We observe that FTL, FTRL, AdaFTRL are the best performing methods, while AdaGrad has worse performance.