

---

# MobILE: Model-Based Imitation Learning From Observation Alone

---

**Rahul Kidambi\***

Amazon Search & AI  
Berkeley CA 94704.  
rk773@cornell.edu

**Jonathan D. Chang**

CS Department, Cornell University  
Ithaca NY 14853.  
jdc396@cornell.edu

**Wen Sun**

CS Department, Cornell University  
Ithaca NY 14853.  
ws455@cornell.edu

## Abstract

This paper studies Imitation Learning from Observations alone (ILFO) where the learner is presented with expert demonstrations that consist only of states visited by an expert (without access to actions taken by the expert). We present a provably efficient model-based framework MobILE to solve the ILFO problem. MobILE involves carefully trading off strategic exploration against imitation - this is achieved by integrating the idea of optimism in the face of uncertainty into the distribution matching imitation learning (IL) framework. We provide a unified analysis for MobILE, and demonstrate that MobILE enjoys strong performance guarantees for classes of MDP dynamics that satisfy certain well studied notions of structural complexity. We also show that the ILFO problem is *strictly harder* than the standard IL problem by presenting an exponential sample complexity separation between IL and ILFO. We complement these theoretical results with experimental simulations on benchmark OpenAI Gym tasks that indicate the efficacy of MobILE. Code for implementing the MobILE framework is available at <https://github.com/rahulkidambi/MobILE-NeurIPS2021>.

## 1 Introduction

This paper considers *Imitation Learning from Observation Alone (ILFO)*. In ILFO, the learner is presented with sequences of states encountered by the expert, *without* access to the actions taken by the expert, meaning approaches based on a reduction to supervised learning (e.g., Behavior cloning (BC) [49], DAgger [50]) are not applicable. ILFO is more general and has potential for applications where the learner and expert have different action spaces, applications like sim-to-real [56, 14] etc.

Recently, [59] reduced the ILFO problem to a sequence of one-step distribution matching problems that results in obtaining a non-stationary policy. This approach, however, is sample inefficient for longer horizon tasks since the algorithm does not effectively reuse previously collected samples when solving the current sub-problem. Another line of work considers model-based methods to infer the expert's actions with either an inverse dynamics [63] or a forward dynamics [16] model; these recovered actions are then fed into an IL approach like BC to output the final policy. These works rely on stronger assumptions that are only satisfied for Markov Decision Processes (MDPs) with injective transition dynamics [68]; we return to this in the related works section.

---

\*Work initiated when RK was a post-doc at Cornell University; work done outside Amazon.

We introduce **MobILE**—**Model-based Imitation Learning and Exploring**, a model-based framework, to solve the ILFO problem. In contrast to existing model-based efforts, MobILE learns the forward transition dynamics model—a quantity that is well defined for any MDP. Importantly, MobILE *combines strategic exploration with imitation* by interleaving a model learning step with a bonus-based, optimistic distribution matching step – a perspective, to the best of our knowledge, that has not been considered in Imitation Learning. MobILE has the ability to automatically trade-off exploration and imitation. It simultaneously explores to collect data to refine the model and imitates the expert wherever the learned model is accurate and certain. At a high level, our theoretical results and experimental studies demonstrate that *systematic exploration is beneficial for solving ILFO reliably and efficiently*, and *optimism* is a both theoretically sound and practically effective approach for strategic exploration in ILFO (see Figure 1 for comparisons with other ILFO algorithms). This paper extends the realm of partial information problems (e.g. Reinforcement Learning and Bandits) where optimism has been shown to be crucial in obtaining strong performance, both in theory (e.g.,  $E^3$  [30], UCB [3]) and practice (e.g., RND [10]). This paper proves that incorporating optimism into the min-max IL framework [69, 22, 59] is *beneficial* for both the theoretical foundations and empirical performance of ILFO.

**Our Contributions:** We present MobILE (Algorithm 1), a provably efficient, model-based framework for ILFO that offers competitive results in benchmark gym tasks. MobILE can be instantiated with various implementation choices owing to its modular design. This paper’s contributions are:

1. The MobILE framework combines ideas of model-based learning, optimism for exploration, and adversarial imitation learning. MobILE achieves global optimality with near-optimal regret bounds for classes of MDP dynamics that satisfy certain well studied notions of complexity. The key idea of MobILE is to use optimism to *trade-off imitation and exploration*.
2. We show an exponential sample complexity gap between ILFO and classic IL where one has access to expert’s actions. This indicates that ILFO is *fundamentally harder* than IL. Our lower bound on ILFO also indicates that to achieve near optimal regret, one needs to perform systematic exploration rather than random or no exploration, both of which will incur sub-optimal regret.
3. We instantiate MobILE with a model ensemble of neural networks and a disagreement-based bonus. We present experimental results on benchmark OpenAI Gym tasks, indicating MobILE compares favorably to or outperforms existing approaches. Ablation studies indicate that optimism indeed helps in significantly improving the performance in practice.

## 1.1 Related Works

**Imitation Learning (IL)** is considered through the lens of two types of approaches: (a) behavior cloning (BC) [45] which casts IL as a reduction to supervised or full-information online learning [49, 50], or, (b) (adversarial) inverse RL [40, 1, 69, 17, 22, 29, 18], which involves minimizing various distribution divergences to solve the IL problem, either with the transition dynamics known (e.g., [69]), or unknown (e.g., [22]). MobILE does not assume knowledge of the transition dynamics, is model-based, and operates without access to the expert’s actions.

**Imitation Learning from Observation Alone (ILFO)** [59] presents a model-free approach FAIL that outputs a non-stationary policy by reducing the ILFO problem into a sequence of min-max problems, one per time-step. While being theoretically sound, this approach cannot share data across different time steps and thus is not data efficient for long horizon problems. Also FAIL in theory only works for discrete actions. In contrast, our paper learns a stationary policy using model-based approaches by reusing data across all time steps and extends to continuous action space. Another line of work [63, 16, 66] relies on learning an estimate of expert action, often through the use of an inverse dynamics models,  $P^e(a|s, s')$ . Unfortunately, an inverse dynamics model is not well defined

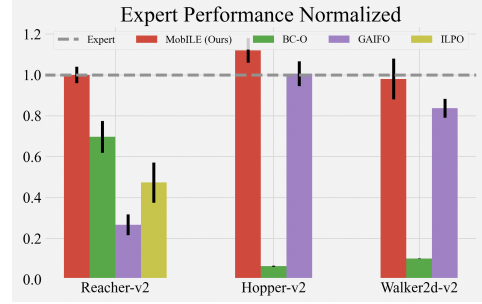


Figure 1: Expert performance normalized scores of ILFO algorithms averaged across 5 seeds in environments with discrete action spaces (Reacher-v2) and continuous action spaces (Hopper-v2 and Walker2d-v2).

in many benign problem instances. For instance, [68, remark 1, section 9.3] presents an example showing that inverse dynamics isn't well defined except in the case when the MDP dynamics is injective (i.e., no two actions could lead to the same next state from the current state. Note that even deterministic transition dynamics doesn't imply injectivity of the MDP dynamics). Furthermore, ILPO [16] applies to MDPs with deterministic transition dynamics and discrete actions. MobILE, on the other hand, learns the forward dynamics model which is always unique and well-defined for both deterministic and stochastic transitions and works with discrete and continuous actions. Another line of work in ILFO revolves around using hand-crafted cost functions that may rely on task-specific knowledge [44, 4, 53]. The performance of policy outputted by these efforts relies on the quality of the engineered cost functions. In contrast, MobILE does not require cost function engineering.

**Model-Based RL** has seen several advances [61, 36, 13] including ones based on deep learning (e.g., [34, 19, 38, 24, 37, 65]). Given MobILE's modularity, these advances in model-based RL can be translated to improved algorithms for the ILFO problem. MobILE bears parallels to provably efficient model-based RL approaches including  $E^3$  [31, 27], R-MAX [7], UCRL [23], UCBVI [5], Linear MDP [67],  $LC^3$  [25], Witness rank [58] which utilize optimism based approaches to trade-off exploration and exploitation. Our work utilizes optimism to trade-off *exploration and imitation*.

## 2 Setting

We consider episodic finite-horizon MDP  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, P^*, H, c, s_0\}$ , where  $\mathcal{S}, \mathcal{A}$  are the state and action space,  $P^* : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  is the MDP's transition kernel,  $H$  is the horizon,  $s_0$  is a fixed initial state (note that our work generalizes when we have a distribution over initial states), and  $c$  is the *state-dependent* cost function  $c : \mathcal{S} \mapsto [0, 1]$ . Our result can be extended to the setting where  $c : \mathcal{S} \times \mathcal{S} \mapsto [0, 1]$ , i.e., the ground truth cost  $c(s, s')$  depends on state and next state pairs. For analysis simplicity, we focus on  $c : \mathcal{S} \mapsto [0, 1]$ .<sup>2</sup>

We denote  $d_P^\pi \in \Delta(\mathcal{S} \times \mathcal{A})$  as the average state-action distribution of policy  $\pi$  under the transition kernel  $P$ , i.e.,  $d_P^\pi(s, a) := \frac{1}{H} \sum_{t=1}^H Pr(s_t = s, a_t = a | s_0, \pi, P)$ , where  $Pr(s_t = s, a_t = a | s_0, \pi, P)$  is the probability of reaching  $(s, a)$  at time step  $t$  starting from  $s_0$  by following  $\pi$  under transition kernel  $P$ . We abuse notation and write  $s \sim d_P^\pi$  to denote a state  $s$  is sampled from the state-wise distribution which marginalizes action over  $d_P^\pi(s, a)$ , i.e.,  $d_P^\pi(s) := \frac{1}{H} \sum_{t=1}^H Pr(s_t = s | s_0, \pi, P)$ . For a given cost function  $f : \mathcal{S} \mapsto [0, 1]$ ,  $V_{P,f}^\pi$  denotes the expected total cost of  $\pi$  under transition  $P$  and cost function  $f$ . Similar to IL setting, in ILFO, the *ground truth cost  $c$  is unknown*. Instead, we can query the expert, denoted as  $\pi^e : \mathcal{S} \mapsto \Delta(\mathcal{A})$ . Note that the expert  $\pi^e$  could be stochastic and does not have to be the optimal policy. The expert, when queried, provides state-only demonstrations  $\tau = \{s_0, s_1 \dots s_H\}$ , where  $s_{t+1} \sim P^*(\cdot | s_t, a_t)$  and  $a_t \sim \pi^e(\cdot | s_t)$ .

The goal is to leverage expert's state-wise demonstrations to learn a policy  $\pi$  that performs as well as  $\pi^e$  in terms of optimizing the ground truth cost  $c$ , with polynomial sample complexity on problem parameters such as horizon, number of expert samples and online samples and underlying MDP's complexity measures (see section 4 for precise examples). We track the progress of any (randomized) algorithm by measuring the (expected) regret incurred by a policy  $\pi$  defined as  $E[V^\pi] - V^{\pi^*}$  as a function of number of online interactions utilized by the algorithm to compute  $\pi$ .

### 2.1 Function Approximation Setup

Since the ground truth cost  $c$  is unknown, we utilize the notion of a function class (i.e., discriminators)  $\mathcal{F} \subset \mathcal{S} \mapsto [0, 1]$  to define the costs that can then be utilized by a planning algorithm (e.g. NPG [26]) for purposes of distribution matching with expert states. If the ground truth  $c$  depends  $(s, s')$ , we use discriminators  $\mathcal{F} \subset \mathcal{S} \times \mathcal{S} \mapsto [0, 1]$ . Furthermore, we use a model class  $\mathcal{P} \subset \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$  to capture the ground truth transition  $P^*$ . For the theoretical results in the paper, we assume realizability:

**Assumption 1.** Assume  $\mathcal{F}$  and  $\mathcal{P}$  captures ground truth cost and transition, i.e.,  $c \in \mathcal{F}$ ,  $P^* \in \mathcal{P}$ .

We will use Integral probability metric (IPM) with  $\mathcal{F}$  as our divergence measure. Note that if  $c \in \mathcal{F}$  and  $c : \mathcal{S} \mapsto [0, 1]$ , then IPM defined as  $\max_{f \in \mathcal{F}} \mathbb{E}_{s \sim d^\pi} f(s) - \mathbb{E}_{s \sim d^{\pi^e}} f(s)$  directly upper

<sup>2</sup>Without any additional assumptions, in ILFO, learning to optimize action-dependent cost  $c(s, a)$  (or  $c(s, a, s')$ ) is **not possible**. For example, if there are two sequences of actions that generate the same sequence of states, without seeing expert's preference over actions, we do not know which actions to commit to.

---

**Algorithm 1** MobILE: The framework of **Model-based Imitation Learning** and **Exploring** for ILFO

---

- 1: **Require:** IPM class  $\mathcal{F}$ , dynamics model class  $\mathcal{P}$ , policy class  $\Pi$ , bonus function class  $\mathcal{B}$ , expert dataset  $\mathcal{D}_e \equiv \{s_i^e\}_{i=1}^N$ .
  - 2: Initialize policy  $\pi_0 \in \Pi$ , replay buffer  $\mathcal{D}_{-1} = \emptyset$ .
  - 3: **for**  $t = 0, \dots, T - 1$  **do**
  - 4:   Execute  $\pi_t$  in true environment  $P^*$  to get samples  $\tau_t = \{s_k, a_k\}_{k=0}^{H-1} \cup s_H$ . Append to replay buffer  $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \tau_t$ .
  - 5:   **Update model and bonus:**  $\hat{P}_{t+1} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  and  $b_{t+1} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$  using buffer  $\mathcal{D}_t$ .
  - 6:   **Optimistic model-based min-max IL:** obtain  $\pi_{t+1}$  by solving equation (1) with  $\hat{P}_{t+1}, b_{t+1}, \mathcal{D}_e$ .
  - 7: **end for**
  - 8: **Return**  $\pi_T$ .
- 

bounds sub-optimality gap  $V^\pi - V^{\pi^e}$ , where  $V^\pi$  is the expected total cost of  $\pi$  under cost function  $c$ . This justifies why minimizing IPM between two state distributions suffices [22, 59]. Similarly, if  $c$  depends on  $s, s'$ , we can simply minimize IPM between two state-next state distributions, i.e.,  $\max_f \mathbb{E}_{s, s' \sim d^\pi} f(s, s') - \mathbb{E}_{s, s' \sim d^{\pi^e}} f(s, s')$  where discriminators now take  $(s, s')$  as input.<sup>3</sup>

To permit generalization, we require  $\mathcal{P}$  to have bounded complexity. For analytical simplicity, we assume  $\mathcal{F}$  is discrete (but exponentially large), and we require the sample complexity of any PAC algorithm to scale polynomially with respect to its complexity  $\ln(|\mathcal{F}|)$ . The  $\ln |\mathcal{F}|$  complexity can be replaced to bounded conventional complexity measures such as Rademacher complexity and covering number for continuous  $\mathcal{F}$  (e.g.,  $\mathcal{F}$  being a Reproducing Kernel Hilbert Space).

### 3 Algorithm

We introduce MobILE (Algorithm 1) for the ILFO problem. MobILE utilizes (a) a function class  $\mathcal{F}$  for Integral Probability Metric (IPM) based distribution matching, (b) a transition dynamics model class  $\mathcal{P}$  for model learning, (c) a bonus parameterization  $\mathcal{B}$  for exploration, (d) a policy class  $\Pi$  for policy optimization. At every iteration, MobILE (in Algorithm 1) performs the following steps:

1. **Dynamics Model Learning:** execute policy in the environment online to obtain state-action-next state  $(s, a, s')$  triples which are appended to the buffer  $\mathcal{D}$ . Fit a transition model  $\hat{P}$  on  $\mathcal{D}$ .
2. **Bonus Design:** design bonus to incentivize exploration where the learnt dynamics model is uncertain, i.e. the bonus  $b(s, a)$  is large at state  $s$  where  $\hat{P}(\cdot|s, a)$  is uncertain in terms of estimating  $P^*(\cdot|s, a)$ , while  $b(s, a)$  is small where  $\hat{P}(\cdot|s, a)$  is certain.
3. **Imitation-Exploration tradeoff:** Given discriminators  $\mathcal{F}$ , model  $\hat{P}$ , bonus  $b$  and expert dataset  $\mathcal{D}_e$ , perform distribution matching by solving the model-based IPM objective with bonus:

$$\pi_{t+1} \leftarrow \arg \min_{\pi \in \Pi} \max_{f \in \mathcal{F}} L(\pi, f; \hat{P}, b, \mathcal{D}_e) := \mathbb{E}_{(s, a) \sim d_{\hat{P}}^\pi} [f(s) - b(s, a)] - \mathbb{E}_{s \sim \mathcal{D}_e} [f(s)], \quad (1)$$

$$\text{where } \mathbb{E}_{s \sim \mathcal{D}_e} f(s) := \sum_{s \in \mathcal{D}_e} f(s) / |\mathcal{D}_e|.$$

Intuitively, the bonus cancels out discriminator's power in parts of the state space where the dynamics model  $\hat{P}$  is not accurate, thus offering freedom for MobILE to explore. We first explain MobILE's components and then discuss MobILE's key property—which is to trade-off *exploration and imitation*.

#### 3.1 Components of MobILE

This section details MobILE's components.

**Dynamics model learning:** For the model fitting step in line 5, we assume that we get a calibrated model in the sense that:  $\|\hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a)\|_1 \leq \sigma_t(s, a), \forall s, a$  for some uncertainty measure  $\sigma_t(s, a)$ , similar to model-based RL works, e.g. [12]. We discuss ways to estimate  $\sigma_t(s, a)$  in the bonus estimation below. There are many examples (discussed in Section 4) that permit efficient

---

<sup>3</sup>we slightly abuse notation here and denote  $d^\pi$  as the average state-next state distribution of  $\pi$ , i.e.,  $d^\pi(s, s') := d^\pi(s) \int_a \pi(a|s) da P^*(s'|s, a)$ .

estimation of these quantities including tabular MDPs, Kernelized nonlinear regulator, nonparametric model such as Gaussian Processes. Consider a general function class  $\mathcal{G} \subset \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$ , one can learn  $\hat{g}_t$  via solving a regression problem, i.e.,

$$\hat{g}_t = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{s, a, s' \in \mathcal{D}_t} \|g(s, a) - s'\|_2^2, \quad (2)$$

and setting  $\hat{P}_t(\cdot|s, a) = \mathcal{N}(\hat{g}_t(s, a), \sigma^2 I)$ , where,  $\sigma$  is the standard deviation of error induced by  $\hat{g}_t$ . In practice, such parameterizations have been employed in several settings in RL with  $\mathcal{G}$  being a multi-layer perceptron (MLP) based function class (e.g., [48]). In Section 4, we also connect this with prior works in provable model-based RL literature.

**Bonus:** We utilize bonuses as a means to incentivize the policy to efficiently explore unknown parts of the state space for improved model learning (and hence better distribution matching). With the uncertainty measure  $\sigma_t(s, a)$  obtained from calibrated model fitting, we can simply set the bonus  $b_t(s, a) = O(H\sigma_t(s, a))$ . How do we obtain  $\sigma_t(s, a)$  in practice? For a general class  $\mathcal{G}$ , given the least square solution  $\hat{g}_t$ , we can define a version space  $\mathcal{G}_t$  as:  $\mathcal{G}_t = \left\{ g \in \mathcal{G} : \sum_{i=0}^{t-1} \sum_{h=0}^{H-1} \|g(s_h^i, a_h^i) - \hat{g}_t(s_h^i, a_h^i)\|_2^2 \leq z_t \right\}$ , with  $z_t$  being a hyper parameter. The version space  $\mathcal{G}_t$  is an *ensemble of functions*  $g \in \mathcal{G}$  which has training error on  $\mathcal{D}_t$  almost as small as the training error of the least square solution  $\hat{g}_t$ . In other words, version space  $\mathcal{G}_t$  contains functions that agree on the training set  $\mathcal{D}_t$ . The uncertainty measure at  $(s, a)$  is then the *maximum disagreement* among models in  $\mathcal{G}_t$ , with  $\sigma_t(s, a) \propto \sup_{f_1, f_2 \in \mathcal{G}_t} \|f_1(s, a) - f_2(s, a)\|_2$ . Since  $g \in \mathcal{G}_t$  agree on  $\mathcal{D}_t$ , a large  $\sigma_t(s, a)$  indicates  $(s, a)$  is novel. See example 3 for more theoretical details.

Empirically, disagreement among an ensemble [41, 6, 11, 43, 37] is used for designing bonuses that incentivize exploration. We utilize an neural network ensemble, where each model is trained on  $\mathcal{D}_t$  (via SGD on squared loss Eq. 2) with different initialization. This approximates the version space  $\mathcal{G}_t$ , and the bonus is set as a function of maximum disagreement among the ensemble’s predictions.

**Optimistic model-based min-max IL:** For model-based imitation (line 6), MobILE takes the current model  $\hat{P}_t$  and the discriminators  $\mathcal{F}$  as inputs and performs policy search to minimize the divergence defined by  $\hat{P}_t$  and  $\mathcal{F}$ :  $d_t(\pi, \pi^e) := \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} (f(s) - b_t(s, a)) - \mathbb{E}_{s \sim d_{\pi^e}} f(s) \right]$ . Note that, for a fixed  $\pi$ , the  $\arg \max_{f \in \mathcal{F}}$  is identical with or without the bonus term, since  $\mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} b_t(s, a)$  is independent of  $f$ . In our implementation, we use the Maximum Mean Discrepancy (MMD) with a Radial Basis Function (RBF) kernel to model discriminators  $\mathcal{F}$ .<sup>4</sup> We compute  $\arg \min_{\pi} d_t(\pi, \pi^e)$  by iteratively (1) computing the  $\arg \max$  discriminator  $f$  given the current  $\pi$ , and (2) using policy gradient methods (e.g., TRPO) to update  $\pi$  inside  $\hat{P}_t$  with  $f - b_t$  as the cost. Specifically, to find  $\pi_t$  (line 6), we iterate between the following two steps:

1. Cost update:  $\hat{f} = \arg \max_{f \in \mathcal{F}} \mathbb{E}_{s \sim d_{\hat{P}_t}^\pi} f(s) - \mathbb{E}_{s \sim \mathcal{D}^e} f(s)$ ,    2. PG Step:  $\hat{\pi} = \hat{\pi} - \eta \cdot \nabla_{\pi} V_{\hat{P}_t, \hat{f} - b_t}^{\hat{\pi}}$ ,

where the PG step uses the learnt dynamics model  $\hat{P}_t$  and the optimistic IPM cost  $\hat{f}(s) - b_t(s, a)$ . Note that for MMD, the cost update step has a closed-form solution.

### 3.2 Exploration And Imitation Tradeoff

We note that MobILE is performing an automatic *trade-off between exploration and imitation*. More specifically, the bonus is designed such that it has high values in the state space that have not been visited, and low values in the state space that have been frequently visited by the sequence of learned policies so far. Thus, by incorporating the bonus into the discriminator  $f \in \mathcal{F}$  (e.g.,  $\hat{f}(s, a) = f(s) - b_t(s, a)$ ), we diminish the power of discriminator  $f$  at novel state-action space regions, which relaxes the state-matching constraint (as the bonus cancels the penalty from the discriminators) at those novel regions so that exploration is encouraged. For well explored states, we force the learner’s states to match the expert’s using the full power of the discriminators. Our work uses optimism (via coupling bonus and discriminators) to carefully balance imitation and exploration.

<sup>4</sup>For MMD with kernel  $k$ ,  $\mathcal{F} = \{w^\top \phi(s, a) \mid \|w\|_2 \leq 1\}$  where  $\phi: \langle \phi(s, a), \phi(s', a') \rangle = k((s, a), (s', a'))$ .



## 4 Analysis

This section presents a general theorem for MobILE that uses the notion of *information gain* [57], and then specializes this result to common classes of stochastic MDPs such as discrete (tabular) MDPs, Kernelized nonlinear regulator [28], and general function class with bounded Eluder dimension [51].

Recall, Algorithm 1 generates one state-action trajectory  $\tau^t := \{s_h^t, a_h^t\}_{h=0}^H$  at iteration  $t$  and estimates model  $\hat{P}_t$  based on  $\mathcal{D}_t = \tau^0, \dots, \tau^{t-1}$ . We present our theorem under the assumption that model fitting gives us a model  $\hat{P}$  and a confidence interval of the model's prediction.

**Assumption 2** (Calibrated Model). *For all iteration  $t$  with  $t \in \mathbb{N}$ , with probability  $1 - \delta$ , we have a model  $\hat{P}_t$  and its associated uncertainty measure  $\sigma_t : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^+$ , such that for all  $s, a \in \mathcal{S} \times \mathcal{A}$* <sup>5</sup>

$$\left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \leq \min \{ \sigma_t(s, a), 2 \}.$$

Assumption 2 has featured in prior works (e.g., [12]) to prove regret bounds in model-based RL. Below we demonstrate examples that satisfy the above assumption.

**Example 1** (Discrete MDPs). *Given  $\mathcal{D}_t$ , denote  $N(s, a)$  as the number of times  $(s, a)$  appears in  $\mathcal{D}_t$ , and  $N(s, a, s')$  number of times  $(s, a, s')$  appears in  $\mathcal{D}_t$ . We can set  $\hat{P}_t(s'|s, a) = N(s, a, s')/N(s, a), \forall s, a, s'$ . We can set  $\sigma_t(s, a) = \tilde{O}\left(\sqrt{S/N(s, a)}\right)$ .*

**Example 2** (KNRs [28]). *For KNR, we have  $P^*(\cdot|s, a) = \mathcal{N}(W^* \phi(s, a), \sigma^2 I)$  where feature mapping  $\phi(s, a) \in \mathbb{R}^d$  and  $\|\phi(s, a)\|_2 \leq 1$  for all  $s, a$ .<sup>6</sup> We can learn  $\hat{P}_t$  via Kernel Ridge regression, i.e.,  $\hat{g}_t(s, a) = \hat{W}_t \phi(s, a)$  where*

$$\hat{W}_t = \operatorname{argmin}_W \sum_{s, a, s' \in \mathcal{D}_t} \|W \phi(s, a) - s'\|_2^2 + \lambda \|W\|_F^2$$

where  $\|\cdot\|_F$  is the Frobenius norm. The uncertainty measure  $\sigma_t(s, a) = \frac{\beta_t}{\sigma} \|\phi(s, a)\|_{\Sigma_t^{-1}}$ ,  $\beta_t = \{2\lambda \|W^*\|_2^2 + 8\sigma^2 \cdot [d_s \ln(5) + 2\ln(t^2/\delta) + \ln(4) + \ln(\det(\Sigma_t)/\det(\lambda I))]\}^{1/2}$ , and  $\Sigma_t = \sum_{k=0}^{t-1} \sum_{h=1}^{H-1} \phi(s_h^k, a_h^k) \phi(s_h^k, a_h^k)^\top + \lambda I$  with  $\lambda > 0$ . See Proposition 12 for more details.

Similar to RKHS, Gaussian processes (GPs) offers a calibrated model [57]. Note that GPs offer similar regret bounds as RKHS; so we do not discuss GPs and instead refer readers to [12].

**Example 3** (General class  $\mathcal{G}$ ). *In this case, assume we have  $P^*(\cdot|s, a) = \mathcal{N}(g^*(s, a), \sigma^2 I)$  with  $g^* \in \mathcal{G}$ . Assume  $\mathcal{G}$  is discrete (but could be exponentially large with complexity measure,  $\ln(|\mathcal{G}|)$ ), and  $\sup_{g \in \mathcal{G}, s, a} \|g(s, a)\|_2 \leq G \in \mathbb{R}^+$ . Suppose model learning step is done by least square:  $\hat{g}_t = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{k=0}^{t-1} \sum_{h=0}^{H-1} \|g(s_h^k, a_h^k) - s_{h+1}^k\|_2^2$ . Compute a version space  $\mathcal{G}_t = \left\{ g \in \mathcal{G} : \sum_{k=0}^{t-1} \sum_{h=0}^{H-1} \|g(s_h^k, a_h^k) - \hat{g}_t(s_h^k, a_h^k)\|_2^2 \leq z_t \right\}$ , where  $z_t = 2\sigma^2 G^2 \ln(2t^2 |\mathcal{G}|/\delta)$  and use this for uncertainty computation. In particular, set uncertainty  $\sigma_t(s, a) = \frac{1}{\sigma} \max_{g_1 \in \mathcal{G}, g_2 \in \mathcal{G}} \|g_1(s, a) - g_2(s, a)\|_2$ , i.e., the maximum disagreement between any two functions in the version space  $\mathcal{G}_t$ . Refer to Proposition 14 for more details.*

The maximum disagreement above motivates our practical implementation where we use an ensemble of neural networks to approximate the version space and use the maximum disagreement among the models' predictions as the bonus. We refer readers to Section 6 for more details.

### 4.1 Regret Bound

We bound regret with the quantity named *Information Gain*  $\mathcal{I}$  (up to some constant scaling factor) [57]:

$$\mathcal{I}_T := \max_{\text{Alg}} \mathbb{E}_{\text{Alg}} \left[ \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min \{ \sigma_t^2(s_h^t, a_h^t), 1 \} \right], \quad (3)$$

<sup>5</sup>the uncertainty measure  $\sigma_t(s, a)$  will depend on the input failure probability  $\delta$ , which we drop here for notational simplicity. When we introduce specific examples, we will be explicit about the dependence on the failure probability  $\delta$  which usually is in the order of  $\ln(1/\delta)$ .

<sup>6</sup>The covariance matrix can be generalized to any PSD matrix with bounded condition number.

where Alg is any adaptive algorithm (thus including [Algorithm 1](#)) that maps from history before iteration  $t$  to some policy  $\pi_t \in \Pi$ . After the main theorem, we give concrete examples for  $\mathcal{I}_T$  where we show that  $\mathcal{I}_T$  has extremely mild growth rate with respect to  $T$  (i.e., logarithmic). Denote  $V^\pi$  as the expected total cost of  $\pi$  under the true cost function  $c$  and the real dynamics  $P^*$ .

**Theorem 3** (Main result). *Assume model learning is calibrated (i.e., [Assumption 2](#) holds for all  $t$ ) and [Assumption 1](#) holds. In [Algorithm 1](#), set bonus  $b_t(s, a) := H \min\{\sigma_t(s, a), 2\}$ . There exists a set of parameters, such that after running [Algorithm 1](#) for  $T$  iterations, we have:*

$$\mathbb{E} \left[ \min_{t \in [0, \dots, T-1]} V^{\pi_t} - V^{\pi^e} \right] \leq O \left( \frac{H^{2.5} \sqrt{\mathcal{I}_T}}{\sqrt{T}} + H \sqrt{\frac{\ln(TH|\mathcal{F}|)}{N}} \right).$$

[Appendix A](#) contains proof of [Theorem 3](#). This theorem indicates that as long as  $\mathcal{I}_T$  grows sublinearly  $o(T)$ , we find a policy that is at least as good as the expert policy when  $T$  and  $N$  approach infinity. For any discrete MDP, KNR [\[28\]](#), Gaussian Processes models [\[57\]](#), and general  $\mathcal{G}$  with bounded Eluder dimension ([\[52, 42\]](#)), we can show that the growth rate of  $\mathcal{I}_T$  with respect to  $T$  is mild.

**Corollary 4** (Discrete MDP). *For discrete MDPs,  $\mathcal{I}_T = \tilde{O}(HS^2A)$  where  $S = |\mathcal{S}|$ ,  $A = |\mathcal{A}|$ . Thus:*

$$\mathbb{E} \left[ \min_{t \in [0, \dots, T-1]} V^{\pi_t} - V^{\pi^e} \right] = \tilde{O} \left( \frac{H^3 S \sqrt{A}}{\sqrt{T}} + H \sqrt{\frac{\ln(|\mathcal{F}|)}{N}} \right).$$

Note that [Corollary 4](#) (proof in [Appendix A.1](#)) hold for any MDPs (not just injective MDPs) and any stochastic expert policy. The dependence on  $A, T$  is tight (see lower bound in [4.2](#)). Now we specialize [Theorem 3](#) to continuous MDPs below.

**Corollary 5** (KNRs (Example 2)). *For simplicity, consider the finite dimension setting  $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ . We can show that  $\mathcal{I}_T = \tilde{O}(Hd + Hdd_s + Hd^2)$  (see [Proposition 13](#) for details), where  $d$  is the dimension of the feature  $\phi(s, a)$  and  $d_s$  is the dimension of the state space. Thus, we have <sup>7</sup>*

$$\mathbb{E} \left[ \min_{t \in [0, \dots, T-1]} V^{\pi_t} - V^{\pi^e} \right] = \tilde{O} \left( \frac{H^3 \sqrt{dd_s + d^2}}{\sqrt{T}} + H \sqrt{\frac{\ln(|\mathcal{F}|)}{N}} \right).$$

**Corollary 6** (General  $\mathcal{G}$  with bounded Eluder dimension (Example 3)). *For general  $\mathcal{G}$ , assume that  $\mathcal{G}$  has Eluder-dimension  $d_E(\epsilon)$  (Definition 3 in [\[42\]](#)). Denote  $d_E = d_E(1/TH)$ . The information gain is upper bounded as  $\mathcal{I}_T = O(Hd_E + d_E \ln(T^3 H |\mathcal{G}|) \ln(TH))$  (see [Proposition 16](#)). Thus,*

$$\mathbb{E} \left[ \min_{t \in [0, \dots, T-1]} V^{\pi_t} - V^{\pi^e} \right] = \tilde{O} \left( \frac{H^3 \sqrt{d_E \ln(TH |\mathcal{G}|)}}{\sqrt{T}} + H \sqrt{\frac{\ln(|\mathcal{F}|)}{N}} \right).$$

Thus as long as  $\mathcal{G}$  has bounded complexity in terms of the Eluder dimension [\[52, 42\]](#), MobILE with the maximum disagreement-based optimism leads to near-optimal guarantees.

## 4.2 Exploration in ILFO and the Exponential Gap between IL and ILFO

To show the benefit of strategic exploration over random exploration in ILFO, we present a *novel* reduction of the ILFO problem to a bandit optimization problem, for which strategic exploration is known to be *necessary* [\[9\]](#) for optimal bounds while random exploration is suboptimal; this reduction indicates that benefit of strategic exploration for solving ILFO efficiently. This reduction also demonstrate that there exists an exponential gap in terms of sample complexity between ILFO and classic IL that has access to expert actions. We leave the details of the reduction framework in [Appendix A.4](#). The reduction allows us to derive the following lower bound for any ILFO algorithm.

**Theorem 7.** *There exists an MDP with number of actions  $A \geq 2$ , such that even with infinitely many expert data, any ILFO algorithm must incur expected commutative regret  $\Omega(\sqrt{AT})$ .*

Specifically we rely on the following reduction where solving ILFO, with even infinite expert data, is at least as hard as solving an MAB problem with the known optimal arm's mean reward which itself

<sup>7</sup>We use  $\tilde{O}$  to suppress log term except the  $\ln(|\mathcal{G}|)$  and  $\ln(|\mathcal{F}|)$  which present the complexity of  $\mathcal{F}$  and  $\mathcal{G}$ .

occurs the same worst case  $\sqrt{AT}$  cumulative regret bound as the one in the classic MAB setting. For MAB, it is known that random exploration such as  $\epsilon$ -greedy will occur suboptimal regret  $O(T^{2/3})$ . Thus to achieve optimal  $\sqrt{T}$  rate, one needs to leverage strategic exploration (e.g., optimism).

Methods such as BC for IL have sample complexity that scales as  $\text{poly} \ln(A)$ , e.g., see [2, Theorem 14.3, Chapter 14] which shows that for tabular MDP, BC learns a policy whose performance is  $O(H^2 \sqrt{S \ln(A)/N})$  away from the expert’s performance (here  $S$  is the number of states in the tabular MDP). Similarly, in interactive IL setting, DAgger [50] can also achieve  $\text{poly} \ln(A)$  dependence in sample complexity. The *exponential gap* in the sample complexity dependence on  $A$  between IL and ILFO formalizes the additional difficulty encountered by learning algorithms in ILFO.

## 5 Practical Instantiation of MobILE

We present a brief practical instantiation MobILE’s components with details in Appendix Section C.

**Dynamics model learning:** We employ Gaussian Dynamics Models parameterized by an MLP [48, 32], i.e.,  $\hat{P}(s, a) := \mathcal{N}(h_\theta(s, a), \sigma^2 I)$ , where,  $h_\theta(s, a) = s + \sigma_{\Delta_s} \cdot \text{MLP}_\theta(s_c, a_c)$ , where,  $\theta$  are MLP’s trainable parameters,  $s_c = (s - \mu_s)/\sigma_s$ ,  $a_c = (a - \mu_a)/\sigma_a$  with  $\mu_s, \mu_a$  (and  $\sigma_s, \sigma_a$ ) being the mean of states, actions (and standard deviation of states and actions) in the replay buffer  $\mathcal{D}$ . Next, for  $(s, a, s') \in \mathcal{D}$ ,  $\Delta_s = s' - s$  and  $\sigma_{\Delta_s}$  is the standard deviation of the state differences  $\Delta_s \in \mathcal{D}$ . We use SGD with momentum [60] for training the parameters  $\theta$  of the MLP.

**Discriminator parameterization:** We utilize MMD as our choice of IPM and define the discriminator as  $f(s) = w^\top \psi(s)$ , where,  $\psi(s)$  are Random Fourier Features [46].

**Bonus parameterization:** We utilize the discrepancy between predictions of a pair of dynamics models  $h_{\theta_1}(s, a)$  and  $h_{\theta_2}(s, a)$  for designing the bonus. Empirically, we found that using more than two models in the ensemble offered little to no improvements. Denote the disagreement at any  $(s, a)$  as  $\delta(s, a) = \|h_{\theta_1}(s, a) - h_{\theta_2}(s, a)\|_2$ , and  $\delta_{\mathcal{D}} = \max_{(s, a) \sim \mathcal{D}} \delta(s, a)$  is the max discrepancy of a replay buffer  $\mathcal{D}$ . We set bonus as  $b(s, a) = \lambda \cdot \min(\delta(s, a)/\delta_{\mathcal{D}}$ , where  $\lambda > 0$  is a tunable parameter.

**PG oracle:** We use TRPO [54] to perform incremental policy optimization inside the learned model.

## 6 Experiments

This section seeks to answer the following questions: (1) How does MobILE compare against other benchmark algorithms? (2) How does optimism impact sample efficiency/final performance? (3) How does increasing the number of expert samples impact the quality of policy outputted by MobILE?

We consider tasks from Open AI Gym [8] simulated with Mujoco [62]: Cartpole-v1, Reacher-v2, Swimmer-v2, Hopper-v2 and Walker2d-v2. We train an expert for each task using TRPO [54] until we obtain an expert policy of average value 460, -10, 38, 3000, 2000 respectively. We setup Swimmer-v2, Hopper-v2, Walker2d-v2 similar to prior model-based RL works [33, 39, 38, 48, 32].

We compare MobILE against the following algorithms: Behavior Cloning (BC), GAIL [22], BC-O [63], ILPO [16] (for environments with discrete actions), GAIFO [64]. Furthermore, recall that BC and GAIL utilize both expert states and actions, information that is not available for ILFO. This makes both BC and GAIL idealistic targets for comparing ILFO methods like MobILE against. As reported by Torabi et al. [63], BC outperforms BC-O in all benchmark results. Moreover, our results indicate MobILE outperforms GAIL and GAIFO in terms of sample efficiency. With reasonable amount of parameter tuning, BC serves as a very strong baseline and nearly solves *deterministic* Mujoco environments. We use code released by the authors for BC-O and ILPO. For GAIL we use an open source implementation [21], and for GAIFO, we modify the GAIL implementation as described by the authors. We present our results through (a) learning curves obtained by averaging the progress of the algorithm across 5 seeds, and, (b) bar plot showing expert normalized scores averaged across 5 seeds using the best performing policy obtained with each seed. Normalized score refers to ratio of policy’s score over the expert score (so that expert has normalized score of 1). For Reacher-v2, since the expert policy has a negative score, we add an constant before normalization. More details can be found in Appendix C.



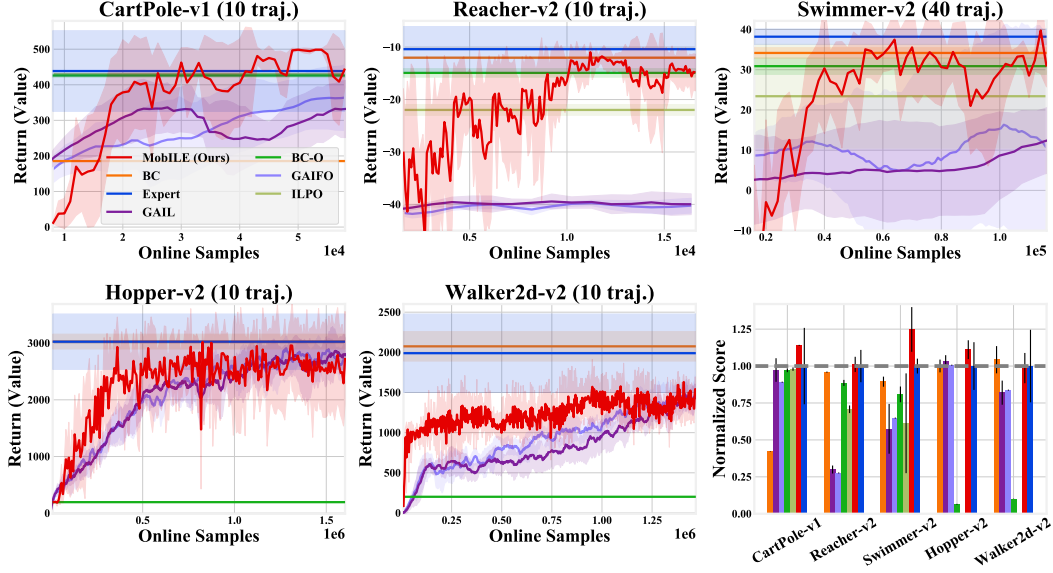


Figure 2: Comparing MobILE (red) against BC (orange), BC-O (green), GAIL (purple), GAIFO (periwinkle), ILPO (green olive). The learning curves are obtained by averaging all algorithms over 5 seeds. MobILE outperforms BC-O, GAIL and matches BC’s behavior despite MobILE not having access to expert actions. The bar plot (bottom-right) presents the best performing policy outputted by each algorithm averaged across 5 seeds for each algorithm. MobILE clearly outperforms BC-O, GAIFO, ILPO while matching the behavior of IL algorithms like BC/GAIL which use expert actions.

### 6.1 Benchmarking MobILE on MuJoCo suite

Figure 2 compares MobILE with BC, BC-O, GAIL, GAIFO and ILPO. MobILE consistently matches or exceeds BC/GAIL’s performance *despite BC/GAIL having access to actions taken by the expert* and MobILE functioning *without* expert action information. MobILE, also, consistently improves upon the behavior of ILFO methods such as BC-O, ILPO, and GAIFO. We see that BC does remarkably well in these benchmarks owing to determinism in the transition dynamics; in the appendix, we consider a variant of the cartpole environment with stochastic dynamics. Our results suggest that BC struggles with stochasticity in the dynamics and fails to solve this task, while MobILE continues to reliably solve this task. Also, note that we utilize 10 expert trajectories for all environments except Swimmer-v2; this is because all algorithms (including MobILE) present results with high variance. We include a learning curve for Swimmer-v2 with 10 expert trajectories in the appendix. The bar plot in Figure 2 shows that within the sample budget shown in the learning curves, MobILE (being a model-based algorithm), presents superior performance in terms of matching expert, thus indicating it is more sample efficient than GAIFO, GAIL (both being model-free methods), ILPO and BC-O.

### 6.2 Importance of the optimistic MDP construction

Figure 3 presents results obtained by running MobILE with and without optimism. In the absence of optimism, the algorithm either tends to be sample inefficient in achieving expert performance or

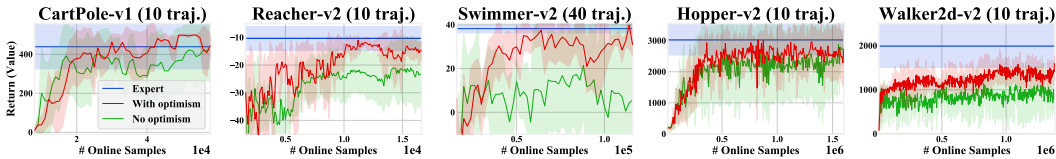


Figure 3: Learning curves obtained by running MobILE with (red) and without (green) optimism. Without optimism, the algorithm learns slowly or does not match the expert, whereas, with optimism, MobILE shows improved behavior by automatically trading off exploration and imitation.

completely fails to solve the problem. Note that without optimism, the algorithm isn’t explicitly incentivized to explore – only implicitly exploring due to noise induced by sampling actions. This, however, is not sufficient to solve the problem efficiently. In contrast, MobILE with optimism presents improved behavior and in most cases, solves the environments with fewer online interactions.

### 6.3 Varying Number of Expert Samples

Table 1 shows the impact of increasing the number of samples drawn from the expert policy for solving the ILFO problem. The main takeaway is that increasing the number of expert samples aids MobILE in reliably solving the problem (i.e. with lesser variance).

Table 1: Expert normalized score and standard deviation of policy outputted by MobILE when varying number of expert trajectories as  $E_1$  and  $E_2$  (specific values represented in parentheses)

Environment	$E_1$	$E_2$	Expert
Cartpole-v1	$1.07 \pm 0.15$ (5)	$1.14 \pm 0$ (10)	$1 \pm 0.25$
Reacher-v2	$1.01 \pm 0.05$ (10)	$0.997 \pm 0.055$ (20)	$1 \pm 0.11$
Swimmer-v2	$1.54 \pm 1.1$ (10)	$1.25 \pm 0.15$ (40)	$1 \pm 0.05$
Hopper-v2	$1.11 \pm 0.064$ (10)	$1.16 \pm 0.03$ (40)	$1 \pm 0.16$
Walker2d-v2	$0.975 \pm 0.12$ (10)	$0.94 \pm 0.038$ (50)	$1 \pm 0.25$

## 7 Conclusions

This paper introduces MobILE, a model-based ILFO approach that is applicable to MDPs with stochastic dynamics and continuous action spaces. MobILE trades-off exploration and imitation, and this perspective is shown to be important for solving the ILFO efficiently both in theory and in practice. Future works include exploring other means for learning dynamics models, performing strategic exploration and extending MobILE to problems with rich observation spaces (e.g. videos).

By not even needing the actions to imitate, ILFO algorithms allow for learning algorithms to capitalize on large amounts of video data available online. Moreover, in ILFO, the learner is successful if it learns to imitate the expert. Any expert policy designed by bad actors can naturally lead to obtaining new policies that continue to imitate and be a negative influence to the society. With this perspective in mind, any expert policy must be thoroughly vetted in order to ensure ILFO algorithms including MobILE are employed in ways that benefit the society.

## Acknowledgements

Rahul Kidambi acknowledges funding from NSF TRIPODS Award CCF – 1740822 at Cornell University. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

## References

- [1] Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*. ACM, 2004.
- [2] Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep*, 2019.
- [3] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [4] Yusuf Aytar, Tobias Pfaff, David Budden, Tom Le Paine, Ziyu Wang, and Nando de Freitas. Playing hard exploration games by watching youtube. In *NeurIPS*, pages 2935–2945, 2018.
- [5] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272, 2017.
- [6] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. In *ITA*, pages 1–9. IEEE, 2018.
- [7] Ronen I. Brafman and Moshe Tennenholtz. R-max - a general polynomial time algorithm for near-optimal reinforcement learning. *J. Mach. Learn. Res.*, 3:213–231, 2001.
- [8] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai Gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [9] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and non-stochastic multi-armed bandit problems. *Found. Trends Mach. Learn.*, 5(1):1–122, 2012.
- [10] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- [11] Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network distillation. In *ICLR*. OpenReview.net, 2019.
- [12] Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient model-based reinforcement learning through optimistic policy search and planning. *arXiv preprint arXiv:2006.08684*, 2020.
- [13] Marc Deisenroth and Carl E. Rasmussen. PILCO: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning*, pages 465–472, 2011.
- [14] Siddharth Desai, Ishan Durugkar, Haresh Karnan, Garrett Warnell, Josiah Hanna, Peter Stone, and AI Sony. An imitation from observation approach to transfer learning with dynamics mismatch. *Advances in Neural Information Processing Systems*, 33, 2020.
- [15] Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- [16] Ashley D. Edwards, Himanshu Sahni, Yannick Schroecker, and Charles L. Isbell Jr. Imitating latent policies from observation. In *ICML*, 2019.
- [17] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *ICML*, 2016.
- [18] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. In *Conference on Robot Learning*, pages 1259–1277. PMLR, 2020.
- [19] Shixiang Gu, Timothy Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous deep q-learning with model-based acceleration, 2016.
- [20] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *CoRR*, abs/1801.01290, 2018.

- [21] 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.
- [22] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *CoRR*, abs/1606.03476, 2016.
- [23] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- [24] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. *CoRR*, abs/1906.08253, 2019.
- [25] Sham Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. *arXiv preprint arXiv:2006.12466*, 2020.
- [26] Sham M. Kakade. A natural policy gradient. In *NIPS*, pages 1531–1538, 2001.
- [27] Sham M. Kakade, Michael J. Kearns, and John Langford. Exploration in metric state spaces. In *ICML*, 2003.
- [28] Sham M. Kakade, Akshay Krishnamurthy, Kendall Lowrey, Motoya Ohnishi, and Wen Sun. Information theoretic regret bounds for online nonlinear control. In *NeurIPS*, 2020.
- [29] Liyiming Ke, Matt Barnes, Wen Sun, Gilwoo Lee, Sanjiban Choudhury, and Siddhartha Srinivasa. Imitation learning as  $f$ -divergence minimization. *arXiv preprint arXiv:1905.12888*, 2019.
- [30] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002.
- [31] Michael Kearns and Satinder Singh. Near optimal reinforcement learning in polynomial time. *Machine Learning*, 49(2-3):209–232, 2002.
- [32] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. *CoRR*, abs/2005.05951, 2020.
- [33] Thanard Kurutach, Ignasi Clavera, Yan Duan, Aviv Tamar, and Pieter Abbeel. Model-ensemble trust-region policy optimization. In *ICLR*. OpenReview.net, 2018.
- [34] Thomas Lampe and Martin A. Riedmiller. Approximate model-assisted neural fitted q-iteration. In *IJCNN*, pages 2698–2704. IEEE, 2014.
- [35] Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- [36] Weiwei Li and Emanuel Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *ICINCO*, pages 222–229, 2004.
- [37] Kendall Lowrey, Aravind Rajeswaran, Sham Kakade, Emanuel Todorov, and Igor Mordatch. Plan Online, Learn Offline: Efficient Learning and Exploration via Model-Based Control. In *International Conference on Learning Representations (ICLR)*, 2019.
- [38] Yuping Luo, Huazhe Xu, Yuezhi Li, Yuandong Tian, Trevor Darrell, and Tengyu Ma. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. *arXiv preprint arXiv:1807.03858*, 2018.
- [39] Anusha Nagabandi, Gregory Kahn, Ronald S. Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *IEEE International Conference on Robotics and Automation*, pages 7559–7566, 2018.
- [40] Andrew Y. Ng and Stuart Russell. Algorithms for inverse reinforcement learning. In *Proc. ICML*, pages 663–670, 2000.

- [41] Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning. *CoRR*, abs/1806.03335, 2018.
- [42] Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the Eluder dimension. In *Advances in Neural Information Processing Systems*, pages 1466–1474, 2014.
- [43] Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement. In *ICML*, pages 5062–5071, 2019.
- [44] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: example-guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graphics*, 2018.
- [45] D. A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. Technical report, CMU, 1989.
- [46] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pages 1177–1184, 2008.
- [47] Aravind Rajeswaran, Kendall Lowrey, Emanuel Todorov, and Sham Kakade. Towards Generalization and Simplicity in Continuous Control. In *NIPS*, 2017.
- [48] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. A game theoretic framework for model based reinforcement learning. *ArXiv*, abs/2004.07804, 2020.
- [49] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In Yee Whye Teh and D. Mike Titterton, editors, *AISTATS*, JMLR Proceedings, pages 661–668, 2010.
- [50] Stéphane Ross, Geoffrey J. Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, pages 627–635, 2011.
- [51] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *NIPS*, pages 2256–2264, 2013.
- [52] Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [53] Karl Schmeckpeper, Oleh Rybkin, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Reinforcement learning with videos: Combining offline observations with interaction. *CoRR*, abs/2011.06507, 2020.
- [54] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR*, abs/1502.05477, 2015.
- [55] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- [56] Yuda Song, Aditi Mavalankar, Wen Sun, and Sicun Gao. Provably efficient model-based policy adaptation. In *International Conference on Machine Learning*, pages 9088–9098. PMLR, 2020.
- [57] Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- [58] Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, pages 2898–2933. PMLR, 2019.
- [59] Wen Sun, Anirudh Vemula, Byron Boots, and Drew Bagnell. Provably efficient imitation learning from observation alone. In *ICML*, volume 97. PMLR, 2019.
- [60] Ilya Sutskever, James Martens, George E. Dahl, and Geoffrey E. Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, volume 28, 2013.



- [61] R. S. Sutton. First results with dyna, an integrated architecture for learning, planning, and reacting. In *Neural Networks for Control*, pages 179–189. The MIT Press: Cambridge, MA, USA, 1990.
- [62] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control. In *IEEE International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.
- [63] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In *IJCAI*, pages 4950–4957, 2018.
- [64] Faraz Torabi, Garrett Warnell, and Peter Stone. Generative adversarial imitation from observation. *arXiv preprint arXiv:1807.06158*, 2018.
- [65] Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- [66] Chao Yang, Xiaojian Ma, Wenbing Huang, Fuchun Sun, Huaping Liu, Junzhou Huang, and Chuang Gan. Imitation learning from observations by minimizing inverse dynamics disagreement. In *NeurIPS*, 2019.
- [67] Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019.
- [68] Zhuangdi Zhu, Kaixiang Lin, Bo Dai, and Jiayu Zhou. Off-policy imitation learning from observations. In *NeurIPS*, 2020.
- [69] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See Section 1
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 4
  - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Section 7
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#) We have read the ethics review guidelines to ensure that this paper conforms to them.
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See section 4
  - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See section 4 and supplemental results for all proofs
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See section 5. All hyperparameters, exact values, and environmental stats are detailed in the supplemental material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See section 6 and supplemental material
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) See section 6. All of our results are presented with error bars across 5 random seeds.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[No\]](#) We do not use any specialized hardware or large scale compute that would be prohibitive for an average end-user to run our experiments. All experiments were run on 1 core and 12 GB of RAM.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 6 for citations for OpenAI Gym and MuJoCo.
  - (b) Did you mention the license of the assets? [\[No\]](#) OpenAI Gym is an opensource software and we paid for our license for MuJoCo.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Related Works . . . . .	2
<b>2</b>	<b>Setting</b>	<b>3</b>
2.1	Function Approximation Setup . . . . .	3
<b>3</b>	<b>Algorithm</b>	<b>4</b>
3.1	Components of MobILE . . . . .	4
3.2	Exploration And Imitation Tradeoff . . . . .	5
<b>4</b>	<b>Analysis</b>	<b>6</b>
4.1	Regret Bound . . . . .	6
4.2	Exploration in ILFO and the Exponential Gap between IL and ILFO . . . . .	7
<b>5</b>	<b>Practical Instantiation of MobILE</b>	<b>8</b>
<b>6</b>	<b>Experiments</b>	<b>8</b>
6.1	Benchmarking MobILE on MuJoCo suite . . . . .	9
6.2	Importance of the optimistic MDP construction . . . . .	9
6.3	Varying Number of Expert Samples . . . . .	10
<b>7</b>	<b>Conclusions</b>	<b>10</b>
<b>A</b>	<b>Analysis of Algorithm 1</b>	<b>17</b>
A.1	Discrete MDPs . . . . .	19
A.2	KNRs . . . . .	20
A.3	General Function Class $\mathcal{G}$ with Bounded Eluder dimension . . . . .	21
A.4	Proof of Theorem 7 . . . . .	23
<b>B</b>	<b>Auxiliary Lemmas</b>	<b>25</b>
<b>C</b>	<b>Implementation Details</b>	<b>25</b>
C.1	Environment Setup and Benchmarks . . . . .	25
C.2	Practical Implementation of MobILE . . . . .	25
C.2.1	Dynamics Model Training . . . . .	26
C.2.2	Replay Buffer . . . . .	26
C.2.3	Design of Bonus Function . . . . .	26
C.2.4	Discriminator Update . . . . .	26
C.2.5	Model-Based Policy Update . . . . .	27
C.3	Hyper-parameter Details . . . . .	27
<b>D</b>	<b>Additional Experimental Results</b>	<b>27</b>

D.1	Modified Cartpole-v0 environment with noise added to transition dynamics . . .	27
D.2	Swimmer Learning Curves . . . . .	28
D.3	Additional Results . . . . .	29
D.4	Ablation Study on Number of Models used for Strategic Exploration Bonus . . . .	29

## A Analysis of Algorithm 1

We start by presenting the proof for the unified main result in [Theorem 3](#). We then discuss the bounds for special instances individually.

The following lemma shows that under [Assumption 2](#), with  $b_t(s, a) = H \min\{\sigma_t(s, a), 2\}$ , we achieve *optimism* at all iterations.

**Lemma 8** (Optimism). *Assume [Assumption 2](#) holds, and set  $b_t(s, a) = H \min\{\sigma_t(s, a), 2\}$ . For all state-wise cost function  $f : \mathcal{S} \mapsto [0, 1]$ , denote the bonus enhance cost as  $\tilde{f}_t(s, a) := f(s) - b_t(s, a)$ . For all policy  $\pi$ , we have the following optimism:*

$$V_{\hat{P}_t, \tilde{f}_t}^\pi \leq V_{P, f}^\pi, \forall t.$$

*Proof.* In the proof, we drop subscript  $t$  for notation simplicity. We consider a fixed function  $f$  and policy  $\pi$ . Also let us denote  $\hat{V}^\pi$  as the value function of  $\pi$  under  $(\hat{P}, \tilde{f})$ , and  $V^\pi$  as the value function under  $(P, f)$ .

Let us start from  $h = H$ , where we have  $\hat{V}_H^\pi(s) = V_H^\pi(s) = 0$ . Assume inductive hypothesis holds at  $h + 1$ , i.e., for any  $s, a$ , we have  $\hat{Q}_{h+1}^\pi(s, a) \leq Q_{h+1}^\pi(s, a)$ . Now let us move to  $h$ . We have:

$$\begin{aligned} \hat{Q}_h^\pi(s, a) - Q_h^\pi(s, a) &= \tilde{f}(s, a) + \mathbb{E}_{s' \sim \hat{P}(\cdot|s, a)} \hat{V}_{h+1}^\pi(s') - f(s) - \mathbb{E}_{s' \sim P(\cdot|s, a)} V_{h+1}^\pi(s') \\ &\leq -H \min\{\sigma(s, a), 2\} + \mathbb{E}_{s' \sim \hat{P}(\cdot|s, a)} V_{h+1}^\pi(s') - \mathbb{E}_{s' \sim P(\cdot|s, a)} V_{h+1}^\pi(s') \\ &\leq -H \min\{\sigma(s, a), 2\} + H \left\| \hat{P}(\cdot|s, a) - P(\cdot|s, a) \right\|_1 \\ &\leq -H \min\{\sigma(s, a), 2\} + H \min\{\sigma(s, a), 2\} = 0, \end{aligned}$$

where the first inequality uses the inductive hypothesis at time step  $h + 1$ . Finally, note that  $V_h^\pi(s) = \mathbb{E}_{a \sim \pi(s)} Q_h^\pi(s, a)$ , which leads to  $\hat{V}_h^\pi(s) \leq V_h^\pi(s)$ . This concludes the induction step.  $\square$

The next lemma concerns the statistical error from finite sample estimation of  $\mathbb{E}_{s \sim d^{\pi^e}} f(s)$ .

**Lemma 9.** *Fix  $\delta \in (0, 1)$ . For all  $t$ , we have that with probability at least  $1 - \delta$ ,*

$$\left| \mathbb{E}_{s \sim d^{\pi^e}} f(s) - \sum_{i=1}^N f(s_i^e) / N \right| \leq 2 \sqrt{\frac{\ln(2t^2 |\mathcal{F}| / \delta)}{N}}, \forall f \in \mathcal{F}.$$

*Proof.* For any  $t$ , we set the failure probability to be  $6\delta/(t^2\pi^2)$  at iteration  $t$  where we abuse notation and point out that  $\pi = 3.14159\dots$ . Thus the total failure probability for all  $t \in \mathbb{N}$  is at most  $\delta$ . We then apply classic Hoeffding inequality to bound  $\mathbb{E}_{s \sim d^{\pi^e}} f(s) - \sum_{i=1}^N f(s_i^e) / N$  with the fact that  $f(s) \in [0, 1]$  for all  $s$ . We conclude the proof by taking a union bound over all  $f \in \mathcal{F}$ .  $\square$

Note that here we have assumed  $s_i^e \sim d^{\pi^e}$  is i.i.d sampled from  $d^{\pi^e}$ . This can easily be achieved by randomly sampling a state from each expert trajectory. Note that we can easily deal with i.i.d trajectories, i.e., if our expert data contains  $N$  many i.i.d trajectories  $\{\tau^1, \dots, \tau^N\}$ , we can apply concentration on the trajectory level, and get:

$$\left| \mathbb{E}_{\tau \sim \pi^e} \left[ \sum_{h=0}^{H-1} f(s_h) \right] - \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} f(s_h^i) \right| \leq O \left( H \sqrt{\frac{\ln(t^2 |\mathcal{F}| / \delta)}{N}} \right),$$

where  $\tau \sim \pi$  denotes that a trajectory  $\tau$  being sampled based on  $\pi$ ,  $s_h^i$  denotes the state at time step  $h$  on the  $i$ -th expert trajectory. Also note that we have  $\mathbb{E}_{s \sim d^\pi} f(s) = \frac{1}{H} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{h=0}^{H-1} f(s_h) \right]$  for any  $\pi, f$ . Together this immediately implies that:

$$\left| \mathbb{E}_{s \sim d^{\pi^e}} f(s) - \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} f(s_h^i) \right| \leq O \left( \sqrt{\frac{\ln(t^2 |\mathcal{F}| / \delta)}{N}} \right),$$

which matches to the bound in [Lemma 9](#).

Now we conclude the proof for [Theorem 3](#).

*Proof of Theorem 3.* Assume that [Assumption 2](#) and the event in [Lemma 9](#) hold. Denote the joint of these two events as  $\mathcal{E}$ . Note that the probability of  $\bar{\mathcal{E}}$  is at most  $2\delta$ . For notation simplicity, denote

$$\epsilon_{stats} = 2\sqrt{\frac{\ln(2T^2 |\mathcal{F}| / \delta)}{N}}.$$

In each model-based planning phase, recall that we perform model-based optimization on the following objective:

$$\pi_t = \operatorname{argmin}_{\pi \in \Pi} \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} [f(s) - b_t(s, a)] - \sum_{i=1}^N f(s_i^e) / N \right].$$

Note that for any  $\pi$ , using the inequality in [Lemma 9](#), we have:

$$\begin{aligned} & \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} (f(s) - b_t(s, a)) - \sum_{i=1}^N f(s_i^e) / N \right] \\ &= \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} (f(s) - b_t(s, a)) - \mathbb{E}_{s \sim d^{\pi^e}} f(s) + \mathbb{E}_{s \sim d^{\pi^e}} f(s) - \sum_{i=1}^N f(s_i^e) / N \right] \\ &\leq \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} (f(s) - b_t(s, a)) - \mathbb{E}_{s \sim d^{\pi^e}} f(s) \right] + \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s \sim d^{\pi^e}} f(s) - \sum_{i=1}^N f(s_i^e) / N \right] \\ &\leq \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^\pi} (f(s) - b_t(s, a)) - \mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi^e}} (f(s) - b_t(s, a)) \right] + \epsilon_{stats} \end{aligned}$$

where in the last inequality we use optimism from [Lemma 8](#), i.e.,  $\mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi^e}} (f(s) - b_t(s, a)) \leq \mathbb{E}_{s \sim d^{\pi^e}} f(s)$ .

Hence, for  $\pi_t$ , since it is the minimizer and  $\pi^e \in \Pi$ , we must have:

$$\begin{aligned} & \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi_t}} (f(s) - b_t(s, a)) - \sum_{i=1}^N f(s_i^e) / N \right] \\ &\leq \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi^e}} (f(s) - b_t(s, a)) - \sum_{i=1}^N f(s_i^e) / N \right] \\ &\leq \max_{f \in \mathcal{F}} \left[ \mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi^e}} (f(s) - b_t(s, a)) - \mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi^e}} (f(s) - b_t(s, a)) \right] + \epsilon_{stats} = \epsilon_{stats}. \end{aligned}$$

Note that  $\mathcal{F}$  contains  $c$ , we must have:

$$\mathbb{E}_{s, a \sim d_{\hat{P}_t}^{\pi_t}} [c(s) - b_t(s, a)] \leq \sum_{i=1}^N c(s_i^e) / N + \epsilon_{stats} \leq \mathbb{E}_{s \sim d^{\pi^e}} c(s) + 2\epsilon_{stats},$$

which means that  $V_{\hat{P}_t; \tilde{c}_t}^{\pi_t} \leq V^{\pi^e} + 2H\epsilon_{stats}$ .

Now we compute the regret in episode  $t$ . First recall that  $b_t(s, a) = H \min\{\sigma_t(s, a), 2\}$ , which means that  $\|b_t\|_\infty \leq 2H$  as  $\|c\|_\infty \leq 1$ , which means that  $\|c - b_t\|_\infty \leq 2H$ . Thus,  $\left\| V_{\hat{P}_t; c - b_t}^\pi \right\|_\infty \leq 2H^2$ .

Recall simulation lemma ([Lemma 18](#)), we have:

$$V^{\pi_t} - V^{\pi^e} \leq V^{\pi_t} - V_{\hat{P}_t; \tilde{c}_t}^{\pi_t} + 2H\epsilon_{stats}$$



$$\begin{aligned}
&= H\mathbb{E}_{s,a \sim d^{\pi_t}} \left[ |\tilde{c}_t(s, a) - c(s)| + 2H^2 \left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \right] + 2H\epsilon_{stat} \\
&= H\mathbb{E}_{s,a \sim d^{\pi_t}} \left[ H \min\{\sigma_t(s, a), 2\} + 2H^2 \left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \right] + 2H\epsilon_{stat} \\
&\leq H\mathbb{E}_{s,a \sim d^{\pi_t}} \left[ H \min\{\sigma_t(s, a), 2\} + 2H^2 \min\{\sigma_t(s, a), 2\} \right] + 2H\epsilon_{stat} \\
&\leq 3H^3\mathbb{E}_{s,a \sim d^{\pi_t}} \min\{\sigma_t(s, a), 2\} + 2H\epsilon_{stat} \\
&\leq 6H^3\mathbb{E}_{s,a \sim d^{\pi_t}} \min\{\sigma_t(s, a), 1\} + 2H\epsilon_{stat}
\end{aligned}$$

Now sum over  $t$ , and denote  $\mathbb{E}_{\pi_t}$  as the conditional expectation conditioned on the history from iteration 0 to  $t - 1$ , we get:

$$\begin{aligned}
\sum_{t=0}^{T-1} [V^{\pi_t} - V^{\pi^e}] &\leq 6H^2 \sum_{t=0}^{T-1} \mathbb{E}_{\pi_t} \left[ \sum_{h=0}^{H-1} \min\{\sigma_t(s_h^t, a_h^t), 1\} \right] + 2HT\epsilon_{stat} \\
&\leq 6H^2 \sum_{t=0}^{T-1} \left[ \sqrt{H} \sqrt{\mathbb{E}_{\pi_t} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\}} \right] + 2HT\epsilon_{stat},
\end{aligned}$$

where in the last inequality we use  $\mathbb{E}[a^\top b] \leq \sqrt{\mathbb{E}[\|a\|_2^2] \mathbb{E}[\|b\|_2^2]}$ .

Recall that  $\pi_t$  are random quantities, add expectation on both sides of the above inequality, and consider the case where  $\mathcal{E}$  holds and  $\bar{\mathcal{E}}$  holds, we have:

$$\begin{aligned}
\mathbb{E} \left[ \sum_{t=0}^{T-1} (V^{\pi_t} - V^{\pi^e}) \right] &\leq 6H^{2.5} \mathbb{E} \left[ \sum_{t=0}^{T-1} \sqrt{\mathbb{E}_{\pi_t} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\}} \right] + 2HT\epsilon_{stat} + \mathbb{P}(\bar{\mathcal{E}})TH \\
&\leq 6H^{2.5} \left[ \sqrt{T} \sqrt{\mathbb{E} \left[ \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\} \right]} \right] + 2HT\epsilon_{stat} + 2\delta TH,
\end{aligned}$$

where in the last inequality, we use  $\mathbb{E}[a^\top b] \leq \sqrt{\mathbb{E}[\|a\|_2^2] \mathbb{E}[\|b\|_2^2]}$ . This implies that that:

$$\mathbb{E} \left[ \min_t V^{\pi_t} - V^{\pi^e} \right] \leq \frac{6H^{2.5}}{\sqrt{T}} \sqrt{\max_{\text{Alg}} \mathbb{E}_{\text{Alg}} \left[ \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\} \right]} + 2H\epsilon_{stat} + 2H\delta.$$

Set  $\delta = 1/(HT)$ , we get:

$$\mathbb{E} [V^\pi - V^{\pi^e}] \leq \frac{6H^{2.5}}{\sqrt{T}} \sqrt{\max_{\text{Alg}} \mathbb{E}_{\text{Alg}} \left[ \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\} \right]} + 2H\sqrt{\frac{\ln(T^3 H |\mathcal{F}|)}{N}} + \frac{2}{T}$$

where Alg is any adaptive mapping that maps from history from  $t = 0$  to the end of the  $t - 1$  iteration to some policy  $\pi_t$ . This concludes the proof.  $\square$

Below we discuss special cases.

## A.1 Discrete MDPs

**Proposition 10** (Discrete MDP Bonus). *With  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ , for all  $t \in \mathbb{N}$ , we have:*

$$\left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \leq \min \left\{ \sqrt{\frac{S \ln(t^2 SA/\delta)}{N_t(s, a)}}, 2 \right\}.$$

*Proof.* The proof simply uses the concentration result for  $\hat{P}_t$  under the  $\ell_1$  norm. For a fixed  $t$  and  $s, a$  pair, using Lemma 6.2 in [2], we have that with probability at least  $1 - \delta$ ,

$$\left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \leq \sqrt{\frac{S \ln(1/\delta)}{N_t(s, a)}}.$$

Applying union bound over all iterations and all  $s, a$  pairs, we conclude the proof.  $\square$

What left is to bound the information gain  $\mathcal{I}$  for the tabular case. For this, we can simply use the [Proposition 13](#) that we develop in the next section for KNR. This is because in KNR, when we set the feature mapping  $\phi(s, a) \in \mathbb{R}^{|S| \cdot |A|}$  to be a one-hot vector with zero everywhere except one in the entry corresponding to  $(s, a)$  pair, the information gain in KNR is reduced to the information gain in the tabular model.

**Proposition 11** (Information Gain in discrete MDPs). *We have:*

$$\mathcal{I}_T = O\left(HS^2A \cdot \ln(TSA/\delta) \ln(1 + TH)\right).$$

*Proof.* Using Lemma B.6 in [\[25\]](#), we have:

$$\sum_{t=0}^{T-1} \min \left\{ \sum_{h=0}^{H-1} \frac{1}{N_t(s_h^t, a_h^t)}, 1 \right\} \leq 2SA \ln(1 + TH).$$

Now using the definition of information gain, we have:

$$\begin{aligned} \mathcal{I}_T &= \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min \left\{ \sigma_t^2(s_h^t, a_h^t), 1 \right\} \leq S \ln(T^2SA/\delta) H \sum_{t=0}^{T-1} \min \left\{ \sum_{h=0}^{H-1} \frac{1}{N_t(s_h^t, a_h^t)}, 1 \right\} \\ &\leq 2HS^2A \ln(T^2SA/\delta) \ln(1 + TH) \end{aligned}$$

This concludes the proof.  $\square$

## A.2 KNRs

Recall the KNR setting from [Example 2](#). The following proposition shows that the bonus designed in [Example 2](#) is valid.

**Proposition 12** (KNR Bonus). *Fix  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ , for all  $t \in \mathbb{N}$ , we have:*

$$\left\| \widehat{P}_t(\cdot | s, a) - P^*(\cdot | s, a) \right\|_1 \leq \min \left\{ \frac{\beta_t}{\sigma} \|\phi(s, a)\|_{\Sigma_t^{-1}}, 2 \right\}, \forall s, a,$$

where  $\beta_t = \sqrt{2\lambda \|W^*\|_2^2 + 8\sigma^2 (d_s \ln(5) + 2 \ln(t^2/\delta) + \ln(4) + \ln(\det(\Sigma_t)/\det(\lambda I)))}$ .

*Proof.* The proof directly follows the confidence ball construction and proof from [\[25\]](#). Specifically, from Lemma B.5 in [\[25\]](#), we have that with probability at least  $1 - \delta$ , for all  $t$ :

$$\left\| (\widehat{W}_t - W^*) (\Sigma_t)^{1/2} \right\|_2^2 \leq \beta_t^2.$$

Thus, with [Lemma 19](#), we have:

$$\left\| \widehat{P}_t(\cdot | s, a) - P^*(\cdot | s, a) \right\|_1 \leq \frac{1}{\sigma} \left\| (\widehat{W}_t - W^*) \phi(s, a) \right\|_2 \leq \left\| (\widehat{W}_t - W^*) (\Sigma_t)^{1/2} \right\| \|\phi(s, a)\|_{\Sigma_t^{-1}} / \sigma \leq \frac{\beta_t}{\sigma} \|\phi(s, a)\|_{\Sigma_t^{-1}}.$$

This concludes the proof.  $\square$

The following proposition bounds the information gain quantity.

**Proposition 13** (Information Gain on KNRs). *For simplicity, let us assume  $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ , i.e.,  $\phi(s, a)$  is a  $d$ -dim feature vector. In this case, we will have:*

$$\mathcal{I}_T = O\left(H \left(d \ln(T^2/\delta) + dd_s + d^2 \ln(1 + \|W^*\|_2^2 TH/\sigma^2)\right) \ln(1 + \|W^*\|_2^2 TH/\sigma^2)\right).$$

*Proof.* From the previous proposition, we know that  $\sigma_t^2(s, a) = (\beta_t^2/\sigma^2) \|\phi(s, a)\|_{\Sigma_t^{-1}}^2$ . Setting  $\lambda = \sigma^2/\|W^*\|_2^2$ , we will have  $\beta_t^2/\sigma^2 \geq 1$ , which means that  $\min\{\sigma_t^2(s, a), 1\} \leq (\beta_t^2/\sigma^2) \min\{\|\phi(s, a)\|_{\Sigma_t^{-1}}^2, 1\}$ .

Note that  $\beta_t$  is non-decreasing with respect to  $t$ , so  $\beta_t \leq \beta_T$  for  $T \geq t$ , where

$$\beta_T = \sqrt{2\sigma^2 + 8\sigma^2(d_s \ln(5) + 2\ln(T^2/\delta) + \ln(4) + d\ln(1 + TH\|W^*\|_2^2/\sigma^2))}$$

Also we have  $\sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min \left\{ \|\phi(s_h^t, a_h^t)\|_{\Sigma_t^{-1}}^2, 1 \right\} \leq H \sum_{t=0}^{T-1} \min \left\{ \sum_{h=0}^{H-1} \|\phi(s_h^t, a_h^t)\|_{\Sigma_t^{-1}}^2, 1 \right\}$ , since  $\min\{a_1, b_1\} + \min\{a_2, b_2\} \leq \min\{a_1 + a_2, b_1 + b_2\}$ . Now call Lemma B.6 in [25], we have:

$$\sum_{t=0}^{T-1} \min \left\{ \sum_{h=0}^{H-1} \|\phi(s_h^t, a_h^t)\|_{\Sigma_t^{-1}}^2, 1 \right\} \leq 2\ln(\det(\Sigma_T)/\det(\lambda I)) = 2d\ln(1 + TH\|W^*\|_2^2/\sigma^2). \quad (4)$$

Finally recall the definition of  $\mathcal{I}_T$ , we have:

$$\begin{aligned} \mathcal{I}_T &= \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min \left\{ \sigma_t^2(s_h^t, a_h^t), 1 \right\} \leq \frac{\beta_T^2}{\sigma^2} \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min \left\{ \|\phi(s_h^t, a_h^t)\|_{\Sigma_t^{-1}}^2, 1 \right\} \leq \frac{\beta_T^2}{\sigma^2} 2Hd\ln(1 + \|W^*\|_2^2 TH/\sigma^2) \\ &\leq 2Hd(2 + 8(d_s \ln(5) + 2\ln(T^2/\delta) + \ln(4) + d\ln(1 + \|W^*\|_2^2 TH/\sigma^2))) \ln(1 + \|W^*\|_2^2 TH/\sigma^2) \\ &= H(4d + 32dd_s + 32d\ln(T^2/\delta) + 32d + 2d^2\ln(1 + \|W^*\|_2^2 TH/\sigma^2)) \ln(1 + \|W^*\|_2^2 TH/\sigma^2), \end{aligned}$$

which concludes the proof.  $\square$

**Extension to Infinite Dimensional RKHS** When  $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{H}$  where  $\mathcal{H}$  is some infinite dimensional RKHS, we can bound our regret using the following intrinsic dimension:

$$\tilde{d} = \max_{\{s_h^t, a_h^t\}_{h=0}^{H-1}\}_{t=0}^{T-1}} \ln \left( I + \frac{1}{\lambda} \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \phi(s_h^t, a_h^t) \phi(s_h^t, a_h^t)^\top \right).$$

In this case, recall Proposition 12, we have:

$$\begin{aligned} \beta_t &\leq \beta_T \leq \sqrt{2\lambda\|W^*\|_2^2 + 8\sigma^2(d_s \ln(5) + 2\ln(t^2/\delta) + \ln(4) + \ln(\det(\Sigma_T)/\det(\lambda I)))} \\ &\leq \sqrt{2\lambda\|W^*\|_2^2 + 8\sigma^2(d_s \ln(5) + 2\ln(t^2/\delta) + \ln(4) + \tilde{d})}. \end{aligned}$$

Also recall Eq. (4), we have:

$$\sum_{t=0}^{T-1} \min \left\{ \sum_{h=0}^{H-1} \|\phi(s_h^t, a_h^t)\|_{\Sigma_t^{-1}}^2, 1 \right\} \leq 2\ln(\det(\Sigma_T)/\det(\lambda I)) \leq 2\tilde{d}.$$

Combine the above two, following similar derivation we had for finite dimensional setting, we have:

$$\mathcal{I}_T = \tilde{O}\left(H\tilde{d}^2 + H\tilde{d}d_s\right).$$

### A.3 General Function Class $\mathcal{G}$ with Bounded Eluder dimension

**Proposition 14.** Fix  $\delta \in (0, 1)$ . Consider a general function class  $\mathcal{G}$  where  $\mathcal{G}$  is discrete, and  $\sup_{g \in \mathcal{G}, s, a} \|g(s, a)\|_2 \leq G$ . At iteration  $t$ , denote  $\hat{g}_t \in \operatorname{argmin}_{g \in \mathcal{G}} \sum_{i=0}^{t-1} \sum_{h=0}^{H-1} \|g(s_h^i, a_h^i) - s_{h+1}^i\|_2^2$ , and denote a version space  $\mathcal{G}_t$  as:

$$\mathcal{G}_t = \left\{ g \in \mathcal{G} : \sum_{i=0}^{t-1} \sum_{h=0}^{H-1} \|g(s_h^i, a_h^i) - \hat{g}_t(s_h^i, a_h^i)\|_2^2 \leq c_t \right\}, \text{ with } c_t = 2\sigma^2 G^2 \ln(2t^2|\mathcal{G}|/\delta).$$

The with probability at least  $1 - \delta$ , we have that for all  $t$ , and all  $s, a$ :

$$\left\| \hat{P}_t(\cdot|s, a) - P^*(\cdot|s, a) \right\|_1 \leq \min \left\{ \frac{1}{\sigma} \max_{g_1 \in \mathcal{G}_t, g_2 \in \mathcal{G}_t} \|g_1(s, a) - g_2(s, a)\|_2, 2 \right\}.$$

*Proof.* Consider a fixed function  $g \in \mathcal{G}$ . Let us denote  $z_h^t = \|g(s_h^t, a_h^t) - s_{h+1}^t\|_2^2 - \|g^*(s_h^t, a_h^t) - s_{h+1}^t\|_2^2$ . We have:

$$\begin{aligned} z_h^t &= (g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t))^\top (g(s_h^t, a_h^t) + g^*(s_h^t, a_h^t) - 2g^*(s_h^t, a_h^t) - 2\epsilon_h^t) \\ &= \|g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t)\|_2^2 - 2(g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t))^\top \epsilon_h^t. \end{aligned}$$

Since  $\epsilon_h^t \sim \mathcal{N}(0, \sigma^2 I)$ , we must have:

$$2(g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t))^\top \epsilon_h^t \sim \mathcal{N}(0, 4\sigma^2 \|g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t)\|_2^2)$$

Since  $\sup_{g, s, a} \|g(s, a)\|_2 \leq G$ , then  $2(g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t))^\top \epsilon_h^t$  is a  $2\sigma G$  sub-Gaussian random variable.

Call Lemma 3 in [52], we have that with probability at least  $1 - \delta$ :

$$\sum_t \sum_h \|g(s_h^t, a_h^t) - s_{h+1}^t\|_2^2 \geq \sum_t \sum_h \|g^*(s_h^t, a_h^t) - s_{h+1}^t\|_2^2 + 2 \sum_t \sum_h \|g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t)\|_2^2 - 4\sigma^2 G^2 \ln(1/\delta).$$

Note that the above can also be derived directly using Azuma-Bernstein's inequality and the property of square loss. With a union bound over all  $g \in \mathcal{G}$ , we have that with probability at least  $1 - \delta$ , for all  $g \in \mathcal{G}$ .

$$\sum_t \sum_h \|g(s_h^t, a_h^t) - s_{h+1}^t\|_2^2 \geq \sum_t \sum_h \|g^*(s_h^t, a_h^t) - s_{h+1}^t\|_2^2 + 2 \sum_t \sum_h \|g(s_h^t, a_h^t) - g^*(s_h^t, a_h^t)\|_2^2 - 4\sigma^2 G^2 \ln(|\mathcal{G}|/\delta).$$

Set  $g = \hat{g}_t$ , and use the fact that  $g_t$  is the minimizer of  $\sum_t \sum_h \|g(s_h^t, a_h^t) - s_{h+1}^t\|_2^2$ , we must have:

$$\sum_t \sum_h \|\hat{g}_t(s_h^t, a_h^t) - g^*(s_h^t, a_h^t)\|_2^2 \leq 2\sigma^2 G^2 \ln(2t^2 |\mathcal{G}|/\delta).$$

Namely we prove that our version space  $\mathcal{G}_t$  contains  $g^*$  for all  $t$ . Thus, we have:

$$\left\| \hat{P}_t(\cdot | s, a) - P^*(\cdot | s, a) \right\|_1 \leq \frac{1}{\sigma} \|\hat{g}_t(s, a) - g^*(s, a)\|_2 \leq \frac{1}{\sigma} \sup_{g_1 \in \mathcal{G}_t, g_2 \in \mathcal{G}_t} \|g_1(s, a) - g_2(s, a)\|_2,$$

where the last inequality holds since both  $g^*$  and  $\hat{g}_t$  belong to the version  $\mathcal{G}_t$ . □

Now we bound the information gain  $\mathcal{I}_T$  below. The proof mainly follows from the proof in [42].

**Lemma 15** (Lemma 1 in [42]). Denote  $\beta_t = 2\sigma^2 G^2 \ln(t^2 |\mathcal{G}|/\delta)$ . Let us denote the uncertainty measure  $w_{t,h} = \sup_{f_1, f_2 \in \mathcal{G}_t} \|f_1(s_h^t, a_h^t) - f_2(s_h^t, a_h^t)\|_2$  (note that  $w_{t,h}$  is non-negative). We have:

$$\sum_{i=0}^{t-1} \sum_{h=0}^{H-1} \mathbf{1}\{w_{t,h}^2 > \epsilon\} \leq \left( \frac{4\beta_t}{\epsilon} + H \right) d_E(\sqrt{\epsilon}).$$

**Proposition 16** (Bounding  $\mathcal{I}_T$ ). Denote  $d = d_E(1/TH)$ . We have

$$\mathcal{I}_T = (1/\sigma^2 + HdG^2/\sigma^2 + 8G^2 \ln(T^2 |\mathcal{G}|/\delta) d \ln(TH)).$$

*Proof.* Note that the uncertainty measures  $w_{t,h}$  are non-negative. Let us reorder the sequence and denote the ordered one as  $w_1 \geq w_2 \geq w_3 \cdots \geq w_{TH-H}$ . For notational simplicity, denote  $M = TH - H$ . We have:

$$\sum_{i=0}^{T-1} \sum_{h=0}^{H-1} w_{t,h}^2 = \sum_{i=0}^{M-1} w_i^2 \leq 1 + \sum_i w_i^2 \mathbf{1}\{w_i^2 \geq \frac{1}{M}\},$$

where the last inequality comes from the fact that  $\sum_i w_i^2 \mathbf{1}\{w_i^2 < 1/M\} \leq M \frac{1}{M} = 1$ . Consider any  $w_t$  where  $w_t^2 \geq 1/M$ . In this case, we know that  $w_1^2 \geq w_2^2 \geq \cdots \geq w_t^2 \geq 1/M$ . This means that:

$$t \leq \sum_i \sum_h \mathbf{1}\{w_{t,h}^2 > w_t^2\} \leq \left( \frac{4\beta_T}{w_t^2} + H \right) d_E(\sqrt{w_t}) \leq \left( \frac{4\beta_T}{w_t^2} + H \right) d_E(1/M),$$

where the second inequality uses the lemma above, and the last inequality uses the fact that  $d_E(\epsilon)$  is non-decreasing when  $\epsilon$  gets smaller. Denote  $d = d_E(1/M)$ . The above inequality indicates that  $w_t^2 \leq \frac{4\beta_T d}{t-Hd}$ . This means that for any  $w_t^2 \geq 1/M$ , we must have  $w_t^2 \leq 4\beta_T d/(t-Hd)$ . Thus, we have:

$$\begin{aligned} \sum_{i=0}^{T-1} \sum_{h=0}^{H-1} w_{t,h}^2 &\leq 1 + HdG^2 + \sum_{\tau=Hd+1}^M w_\tau^2 \mathbf{1}\{w_\tau^2 \geq 1/M\} \leq 1 + HdG^2 + 4\beta_T d \ln(M) \\ &= 1 + HdG^2 + 4\beta_T d \ln(TH). \end{aligned}$$

Finally, recall the definition of  $\mathcal{I}_T$ , we have:

$$\sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \min\{\sigma_t^2(s_h^t, a_h^t), 1\} \leq \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} \sigma_t^2(s_h^t, a_h^t) \leq \frac{1}{\sigma^2} \sum_{t=0}^{T-1} \sum_{h=0}^{H-1} w_{t,h}^2 \leq \frac{1}{\sigma^2} (1 + HdG^2 + 4\beta_T d \ln(TH)).$$

This concludes the proof.  $\square$

#### A.4 Proof of Theorem 7

This section provides the proof of Theorem 7.

First we present the reduction from a bandit optimization problem to ILFO.

Consider a Multi-armed bandit (MAB) problem with  $A$  many actions  $\{a_i\}_{i=1}^A$ . Each action's ground truth reward  $r_i$  is sampled from a Gaussian with mean  $\mu_i$  and variance 1. Without loss of generality, assume  $a_1$  is the optimal arm, i.e.,  $\mu_1 \geq \mu_i \forall i \neq 1$ . We convert this MAB instance into an MDP. Specifically, set  $H = 2$ . Suppose we have a fixed initial state  $s_0$  which has  $A$  many actions. For the one step transition, we have  $P(\cdot|s_0, a_i) = \mathcal{N}(\mu_i, 1)$ , i.e.,  $g^*(s_0, a_i) = \mu_i$ . Here we denote the optimal expert policy  $\pi^e$  as  $\pi^e(s_0) = a_1$ , i.e., expert policy picks the optimal arm in the MAB instance. Hence, when executing  $\pi^e$ , we note that the state  $s_1$  generated from  $\pi^e$  is simply the stochastic reward of  $a_1$  in the original MAB instance. Assume that we have observed infinitely many such  $s_1$  from the expert policy  $\pi^e$ , i.e., we have infinitely many samples of expert state data, i.e.,  $N \rightarrow \infty$ . Note, however, we do not have the actions taken by the expert (since this is the ILFO setting). This expert data is equivalent to revealing the optimal arm's mean reward  $\mu_1$  to the MAB learner a priori. Hence solving the ILFO problem on this MDP is no easier than solving the original MAB instance with additional information which is that optimal arm's mean reward is  $\mu_1$  (but the best arm's identity is unknown).

Below we show the lower bound for solving the MAB problem where the optimal arm's mean is known.

**Theorem 17.** *Consider best arm identification of Gaussian MAB with the additional information that the optimal arm's mean reward is  $\mu$ . For any algorithm, there exists a MAB instance with number of arms  $A \geq 2$ , such that the expected cumulative regret is still  $\Omega(\sqrt{AT})$ , i.e., the additional information does not help improving the worst-case regret bound to solve the MAB instance.*

*Proof of Theorem 17.* Below, we will construct  $A$  many MAB instances where each instance has  $A$  many arms and each arm has a Gaussian reward distribution with the fixed variance  $\sigma^2$ . Each of the  $A$  instances has the maximum mean reward equal to  $\Delta$ , i.e., all these  $A$  instances have the same maximum arm mean reward. Consider any algorithm Alg that maps  $\Delta$  together with the history of the interactions  $\mathcal{H}_t = \{a_0, r_0, a_1, r_1, \dots, a_{t-1}, r_{t-1}\}$  to a distribution over  $A$  actions. We will show for any such algorithm alg that knows  $\Delta$ , with constant probability, there must exist a MAB instance from the  $A$  many MAB instances, such that Alg suffers at least  $\Omega(\sqrt{AT})$  regret where  $T$  is the number of iterations.

Now we construct the  $A$  instances as follows. Consider the  $i$ -th instance ( $i = 1, \dots, A$ ). For arm  $j$  in the  $i$ -th instance, we define its mean as  $\mu_j^i = \mathbf{1}\{i = j\}\Delta$ . Namely, for MAB instance  $i$ , its arms have mean reward zero everywhere except that the  $i$ -th arm has reward mean  $\Delta$ . Note that all these MAB instances have the same maximum mean reward, i.e.,  $\Delta$ . Hence, we cannot distinguish them by just revealing  $\Delta$  to the learner.

We will construct an additional MAB instance (we name it as 0-th MAB instance) whose arms have reward mean zero. Note that this MAB instance has maximum mean reward 0 which is different from



the previous  $A$  MAB instances that we constructed. However, we will only look at the regret of Alg on the previously constructed  $A$  MAB instances. I.e., we do not care about the regret of Alg( $\Delta, \mathcal{H}_t$ ) on the 0-th MAB instance.

Let us denote  $\mathbb{P}_i$  (for  $i = 0, \dots, A$ ) as the distribution of the outcomes of algorithm Alg( $\Delta, \mathcal{H}_t$ ) interacting with MAB instance  $i$  for  $n$  iterations, and  $\mathbb{E}_j[N_i(T)]$  as the expected number of times arm  $i$  is pulled by Alg( $\Delta, \mathcal{H}_t$ ) in MAB instance  $j$ . Consider MAB instance  $i$  with  $i \geq 1$ :

$$\mathbb{E}_i[N_i(T)] - \mathbb{E}_0[N_i(T)] \leq T \|\mathbb{P}_i - \mathbb{P}_0\|_1 \leq T \sqrt{\text{KL}(\mathbb{P}_0, \mathbb{P}_i)} \leq T \sqrt{\Delta^2 \mathbb{E}_0[N_i(T)]},$$

where the last step uses the fact that we are running the same algorithm Alg( $\Delta, \mathcal{H}_t$ ) on both instance 0 and instance  $i$  (i.e., same policy for generating actions), and thus,  $\text{KL}(\mathbb{P}_0, \mathbb{P}_i) = \sum_{j=1}^A \mathbb{E}_0[N_j(T)] \text{KL}(q_0(j), q_i(j))$  (Lemma 15.1 in [35]), where  $q_i(j)$  is the reward distribution of arm  $j$  at instance  $i$ . Also recall that for instance 0 and instance  $i$ , their rewards only differ at arm  $i$ .

This implies that:

$$\mathbb{E}_i[N_i(T)] \leq \mathbb{E}_0[N_i(T)] + T \sqrt{\Delta^2 \mathbb{E}_0[N_i(T)]}.$$

Sum over  $i = 1, \dots, A$  on both sides, we have:

$$\begin{aligned} \sum_{i=1}^A \mathbb{E}_i[N_i(T)] &\leq T + T \sum_{i=1}^A \sqrt{\Delta^2 \mathbb{E}_0[N_i(T)]} \leq T + T \sqrt{A} \sqrt{\sum_{i=1}^A \Delta^2 \mathbb{E}_0[N_i(T)]} \\ &\leq T + T \sqrt{A} \sqrt{\Delta^2 T} \end{aligned}$$

Now let us calculate the regret of Alg( $\Delta, \mathcal{H}_t$ ) on  $i$ -th instance, we have:

$$R_i = T\Delta - \mathbb{E}_i[N_i(T)]\Delta.$$

Sum over  $i = 1, \dots, A$ , we have:

$$\sum_{i=1}^A R_i = \Delta \left( AT - \sum_{i=1}^A \mathbb{E}_i[N_i(T)] \right) \geq \Delta \left( AT - T - T \sqrt{A \Delta^2 T} \right)$$

Set  $\Delta = c\sqrt{A/T}$  for some  $c$  that we will specify later, we get:

$$\sum_{i=1}^A R_i \geq c \sqrt{\frac{A}{T}} (AT - T - cAT).$$

Set  $c = 1/4$ , we get:

$$\sum_{i=1}^A R_i \geq c \sqrt{\frac{A}{T}} (AT - T - cAT) \geq \frac{1}{4} \sqrt{AT} (A - 1 - A/4) = \frac{1}{4} \sqrt{AT} (3A/4 - 1) \geq \frac{1}{4} \sqrt{AT} (A/4),$$

assuming  $A \geq 2$ .

Thus there must exist  $i \in \{1, \dots, A\}$ , such that:

$$R_i \geq \frac{1}{16} \sqrt{AT}.$$

Note that the above construction considered any algorithm Alg( $\Delta, \mathcal{H}_t$ ) that maps  $\Delta$  and history to action distributions. Thus it concludes the proof.  $\square$

The hardness result in Theorem 17 and the reduction from MAB to ILFO together implies the lower bound for ILFO in Theorem 7, namely solving ILFO with cumulative regret smaller than  $O(\sqrt{AT})$  will contradict the MAB lower bound in Theorem 17.

## B Auxiliary Lemmas

**Lemma 18** (Simulation Lemma). *Consider any two functions  $f : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$  and  $\hat{f} : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ , any two transitions  $P$  and  $\hat{P}$ , and any policy  $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$ . We have:*

$$\begin{aligned} V_{P,f}^\pi - V_{\hat{P},\hat{f}}^\pi &= \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_P^\pi} \left[ f(s,a) - \hat{f}(s,a) + \mathbb{E}_{s' \sim P(\cdot|s,a)} V_{\hat{P},\hat{f};h}^\pi(s') - \mathbb{E}_{s' \sim \hat{P}(\cdot|s,a)} V_{\hat{P},\hat{f};h}^\pi(s') \right] \\ &\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_P^\pi} \left[ f(s,a) - \hat{f}(s,a) + \|V_{\hat{P},\hat{f};h}^\pi\|_\infty \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \right]. \end{aligned}$$

where  $V_{P,f;h}^\pi$  denotes the value function at time step  $h$ , under  $\pi$ ,  $P$ ,  $f$ .

Such simulation lemma is standard in model-based RL literature and can be found, for instance, in the proof of Lemma 10 from [58].

**Lemma 19.** *Consider two Gaussian distribution  $P_1 := \mathcal{N}(\mu_1, \sigma^2 I)$  and  $P_2 := \mathcal{N}(\mu_2, \sigma^2 I)$ . We have:*

$$\|P_1 - P_2\|_1 \leq \frac{1}{\sigma} \|\mu_1 - \mu_2\|_2.$$

The above lemma can be proved by Pinsker’s inequality and the closed-form of the KL divergence between  $P_1$  and  $P_2$ .

## C Implementation Details

### C.1 Environment Setup and Benchmarks

This section sketches the details of how we setup the environments. We utilize the standard environment horizon of 500, 50, 200 for Cartpole-v1, Reacher-v2, Cartpole-v0. For Swimmer-v2, Hopper-v2 and Walker2d-v2, we work with the environment horizon set to 400 [33, 39, 38, 48, 32]. Furthermore, for Hopper-v2, Walker2d-v2, we add the velocity of the center of mass to the state parameterization [48, 38, 32]. As noted in the main text, the expert policy is trained using NPG/TRPO [26, 54] until it hits a value of (approximately) 460, −10, 38, 3000, 2000, 170 for Cartpole-v1, Reacher-v2, Swimmer-v2, Hopper-v2, Walker2d-v2, Cartpole-v0 respectively. Furthermore, for Walker2d-v2 we utilized pairs of states  $(s, s')$  for defining the feature representation used for parameterizing the discriminator. All the results presented in the experiments section are averaged over five seeds. Furthermore, in terms of baselines, we compare MobILE to BC, BC-O, ILPO, GAIL and GAIFO. Note that BC/GAIL has access to expert actions whereas our algorithm does not have access to the expert actions. We report the average of the best performance offered by BC/BC-O when run with five seeds, even if this occurs at different epochs for each of the runs - this gives an upper hand to BC/BC-O. Moreover, note that for BC, we run the supervised learning algorithm for 500 passes. Furthermore, we run BC-O/GAIL with same number of online samples as MobILE in order to present our results. Furthermore, we used 2 CPUs with 16-32 GB of RAM usage to perform all our benchmarking runs implemented in Pytorch. Finally, our codebase utilizes Open-AI’s implementation of TRPO [15] for environments with discrete actions, and the MJRL repository [47] for working with continuous action environments. With regards to results in the main paper, our bar graph presenting normalized results was obtained by dividing every algorithm’s performance (mean/standard deviation) by the expert mean; for Reacher-v2 because the rewards themselves are negative, we first added a constant offset to make all the algorithm’s performance to become positive, then, divided by the mean of expert policy.

### C.2 Practical Implementation of MobILE

We will begin with presenting the implementation details of MobILE (refer to Algorithm 2):



$$\begin{aligned} &\equiv \max_w L_\zeta(w; \pi, \hat{P}, b, \mathcal{D}_e) = \mathbb{E}_{(s,a) \sim d_{\hat{P}}^\pi} [f_w(s) - b(s, a)] - \mathbb{E}_{s \sim \mathcal{D}_e} [f_w(s)] - \frac{1}{2} \cdot (\|w\|_2^2 - \zeta), \\ &\implies \partial_w L_\zeta(w; \pi, \hat{P}, b, \mathcal{D}_e) = \mathbb{E}_{s \sim d_{\hat{P}}^\pi} [\psi(s)] - \mathbb{E}_{s \sim \mathcal{D}_e} [\psi(s)] - w \in 0, \end{aligned}$$

where,  $\partial_w L_\zeta(w; \pi, \hat{P}, b, \mathcal{D}_e)$  denotes the sub-differential of  $L_\zeta(\cdot)$  wrt  $w$ . This in particular implies the following:

1. **Exact Update:**  $w^* = \mathcal{P}_{\mathcal{B}(\zeta)} \left( \mathbb{E}_{s \sim d_{\hat{P}}^\pi} [\psi(s)] - \mathbb{E}_{s \sim \mathcal{D}_e} [\psi(s)] \right)$ ,  $\mathcal{P}$ . is the projection operator, and  $\mathcal{B}(\zeta)$  is the  $\zeta$ -norm ball.
2. **Gradient Ascent Update:**  $w_{t+1} = \mathcal{P}_{\mathcal{B}(\zeta)} \left( (1 - \eta_w)w_t + \eta_w \cdot \left( \mathbb{E}_{s \sim d_{\hat{P}}^\pi} [\psi(s)] - \mathbb{E}_{s \sim \mathcal{D}_e} [\psi(s)] \right) \right)$ ,  $\eta_w > 0$  is the step-size.

We found empirically either of the updates to work reasonably well. In the `Swimmer-v2` task, we use the gradient ascent update with  $\eta_w = 0.67$ , and, in the other tasks, we utilize the exact update. Furthermore, we empirically observe the gradient ascent update to yield more stability compared to the exact updates. In the case of `Walker2d-v2`, we found it useful to parameterize the discriminator based on pairs of states  $(s, s')$ .

### C.2.5 Model-Based Policy Update

Once the maximization of the discriminator parameters  $w$  is performed, consider the policy optimization problem, i.e.,

$$\begin{aligned} \min_{\pi} L(\pi; w, \hat{P}, b, \mathcal{D}_e) &:= \mathbb{E}_{(s,a) \sim d_{\hat{P}}^\pi} [f_w(s) - b(s, a)] - \mathbb{E}_{s \sim \mathcal{D}_e} [f_w(s)] \\ \min_{\pi} L(\pi; w, \hat{P}, b, \mathcal{D}_e) &= \mathbb{E}_{(s,a) \sim d_{\hat{P}}^\pi} [f_w(s) - b(s, a)] \end{aligned}$$

Hence we perform model-based policy optimization under  $\hat{P}$  and cost function  $f_w(s) - b(s, a)$ . In practice, we perform approximate minimization of  $L(\cdot)$  by incrementally updating the policy using  $K_{PG}$ -steps of policy gradient, where,  $K_{PG}$  is a tunable hyper-parameter. In our experiments, we find that setting  $K_{PG}$  to be around 10 to generally be a reasonable choice (for precise values, refer to Table 2). This paper utilizes TRPO [54] as our choice of policy gradient method; note that this can be replaced by other alternatives including PPO [55], SAC [20] *etc.* Similar to practical implementations of existing policy gradient methods, we implement a reward filter by clipping the IPM reward  $f(s)$  by truncating it between  $c_{\min}$  and  $c_{\max}$  as this leads to stability of the policy gradient updates. Note that the minimization is done with access to  $\hat{P}$ , which implies we perform *model-based* planning. Empirically, for purposes of tuning the exploration-imitation parameter  $\lambda$ , we minimize a surrogate namely:  $\mathbb{E}_{(s,a) \sim d_{\hat{P}}^\pi} [(1 - \lambda) \cdot f_w(s) - b(s, a)]$  (recall that  $b(s, a)$  has a factor of  $\lambda$  associated with it). This ensures that we can precisely control the magnitude of the bonuses against the IPM costs, which, in our experience is empirically easier to work with.

### C.3 Hyper-parameter Details

This section presents an overview of the list of hyper-parameters necessary to implement Algorithm 1 in practice, as described in Algorithm 2. The list of hyper-parameters is precisely listed out in Table 2. The hyper-parameters are broadly categorized into ones corresponding to various components of MobILE, namely, (a) environment specifications, (b) dynamics model, (c) ensemble based bonus, (d) IPM parameterization, (e) Policy parameterization, (f) Planning algorithm parameters, (g) Critic parameterization. Note that if there a hyper-parameter that has not been listed, for instance, say, the value of momentum for the ADAM optimizer in the critic, this has been left as is the default value defined in Pytorch.

## D Additional Experimental Results

### D.1 Modified Cartpole-v0 environment with noise added to transition dynamics

Parameter	Cartpole-v1	Reacher-v2	Swimmer-v2	Cartpole-v0	Hopper-v2	Walker2d-v2
<b>Environment Specifications</b>						
Horizon $H$	500	50	400	200	400	400
Expert Performance ( $\approx$ )	460	-10	38	181	3000	2000
# online samples per outer loop	$2 \cdot H$	$2 \cdot H$	$2 \cdot H$	$2 \cdot H$	$8 \cdot H$	$3 \cdot H$
<b>Dynamics Model</b>						
Architecture/Non-linearity	MLP(64, 64)/ReLU	MLP(64, 64)/ReLU	MLP(512, 512)/ReLU	MLP(64, 64)/ReLU	MLP(512, 512)/ReLU	MLP(512, 512)/ReLU
Optimizer(LR, Momentum, Batch Size)	SGD(0.005, 0.99, 256)	SGD(0.005, 0.99, 256)	SGD(0.005, 0.99, 256)	SGD(0.005, 0.99, 256)	SGD(0.005, 0.99, 256)	SGD(0.005, 0.99, 256)
# train passes per outer loop	20	100	100	20	50	200
Grad Clipping	2.0	2.0	1.0	2.0	4.0	1.0
Replay Buffer Size	$10 \cdot H$	$10 \cdot H$	$10 \cdot H$	$10 \cdot H$	$16 \cdot H$	$15 \cdot H$
<b>Ensemble based bonus</b>						
# models/bonus range	$2/[0, 1]$	$2/[0, 1]$	$2/[0, 1]$	$2/[0, 1]$	$2/[0, 1]$	$2/[0, 1]$
<b>IPM parameters</b>						
Step size for $w$ update ( $\eta_w$ )	Exact	Exact	0.33	Exact	Exact	Exact
# RFFs/BW Heuristic	128/0.1 quantile	128 / 0.1 quantile	128 / 0.1 quantile	128 / 0.1 quantile	128 / 0.1 quantile	128 / 0.1 quantile
<b>Policy parameterization</b>						
Architecture/Non-linearity	MLP(64, 64)/TanH	MLP(64, 64)/TanH	MLP(64, 64)/TanH	MLP(32, 32)/TanH	MLP(32, 32)/TanH	MLP(32, 32)/TanH
Policy Constraints	None	None	None	None	$\log \sigma_{\min} = -1.0$	$\log \sigma_{\min} = -2.0$
<b>Planning Algorithm</b>						
# model samples per TRPO step	$2 \cdot H$	$10 \cdot H$	$4 \cdot H$	$4 \cdot H$	$8 \cdot H$	$20 \cdot H$
# TRPO steps per outer loop ( $K_{PG}$ )	3	10	20	5	10	15
TRPO Parameters	(50, 0.001, 0.01, 0.97, 0.995)	(100, 0.001, 0.01, 0.97, 0.995)	(100, 0.001, 0.01, 0.97, 0.995)	(100, 0.001, 0.01, 0.97, 0.995)	(10, 0.0001, 0.025, 0.97, 0.995)	(10, 0.0001, 0.025, 0.97, 0.995)
<b>Critic parameterization</b>						
Architecture/Non-linearity	MLP(128, 128)/ReLU	MLP(128, 128)/ReLU	MLP(128, 128)/ReLU	MLP(32, 32)/ReLU	MLP(128, 128)/ReLU	MLP(128, 128)/ReLU
Optimizer (LR, Batch Size, $\epsilon$ , Regularization)	Adam(0.001, 64, $1e-5$ , 0)	Adam(0.001, 64, $1e-5$ , 0)	Adam(0.001, 64, $1e-5$ , 0)	Adam(0.001, 64, $1e-5$ , 0)	Adam(0.001, 64, $1e-8$ , $1e-3$ )	Adam(0.001, 64, $1e-8$ , $1e-3$ )
# train passes per TRPO update	1	1	1	1	2	2

Table 2: List of various Hyper-parameters employed in MobILE’s implementation.

We consider a stochastic variant of Cartpole-v0, wherein, we add additive Gaussian noise of variance unknown to the learner in order to make the transition dynamics of the environment to be stochastic. Specifically, we train an expert of value  $\approx 170$  in Cartpole-v0 with stochastic dynamics using TRPO. Now, using 20 trajectories drawn from this expert, we wish to consider solving the ILFO problem using MobILE as well as other baselines including BC, BC-O, ILPO, GAIL and GAIFO. Figure 4 presents the result of this comparison. Note that MobILE compares favorably against other baseline methods - in particular, BC tends suffer in environments like Cartpole-v0 with stochastic dynamics because of increased generalization error of the supervised learning algorithm used for learning a policy. Our algorithm is competitive with both BC-O, GAIL, GAIFO and ILPO. Note that BC-O tends to outperform BC both in Cartpole-v1 and in Cartpole-v0 (with stochastic dynamics).

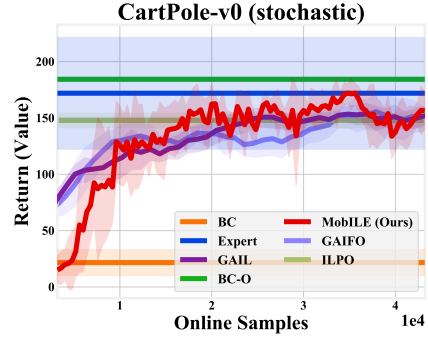


Figure 4: Learning curves for Cartpole-v0 with stochastic dynamics with 20 expert trajectories comparing MobILE with BC, BC-O, GAIL, GAIFO and ILPO.

## D.2 Swimmer Learning Curves

We supplement the learning curves for Swimmer-v2 (with 40 expert trajectories) with the learning curves for Swimmer-v2 with 10 expert trajectories in figure 5. As can be seen, MobILE outperforms baseline algorithms such as BC, BC-O, ILPO, GAIL and GAIFO in Swimmer-v2 with both 40 and 10 expert trajectories. The caveat is that for 10 expert trajectories, all algorithms tend to show a lot more variance in their behavior and this reduces as we move to the 40 expert trajectory case.

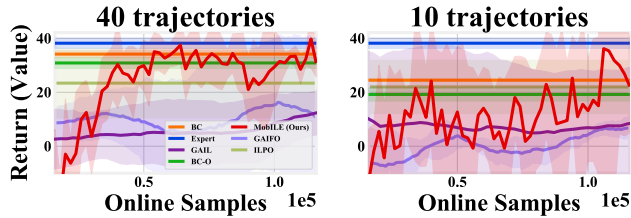


Figure 5: Learning curves for Swimmer-v2 with 40 (left) and 10 (right) expert trajectories comparing MobILE with BC, BC-O, ILPO, GAIL and GAIFO. MobILE continues to perform well relative to all other benchmarks with both 10 and 40 expert trajectories. The variance of the algorithm as well as the benchmarks is notably higher with lesser number of expert trajectories.



### D.3 Additional Results

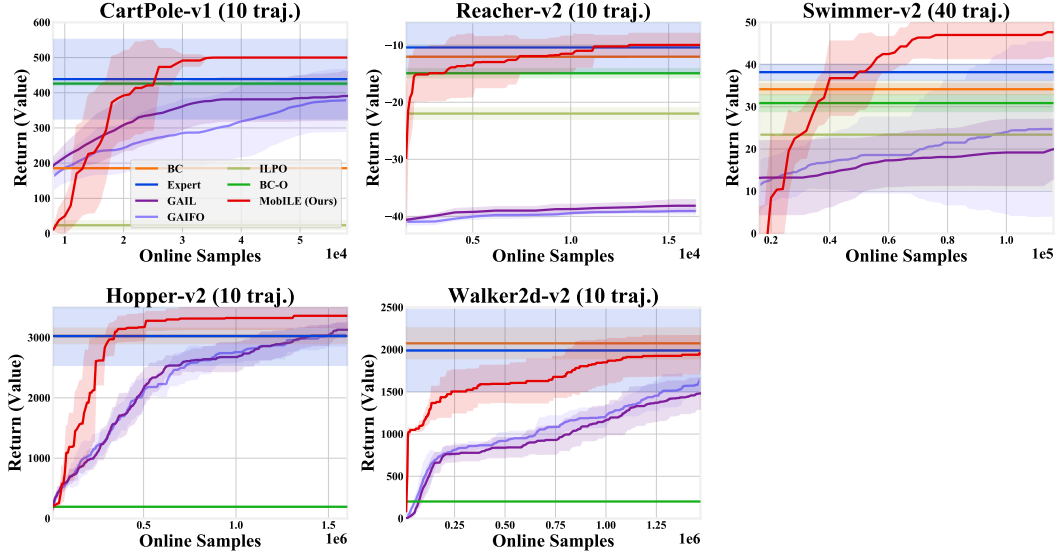


Figure 6: Learning curves tracking the running maximum averaged across seeds comparing MobILE against BC, BC-O, ILPO, GAIL and GAIFO. MobILE tends to reach expert performance consistently and in a more sample efficient manner.

In this section, we give another view of our results for MobILE compared against the baselines (BC/BC-O/ILPO/GAIL/GAIFO) by tracking the running maximum of each policy’s value averaged across seeds. Specifically, for every iteration  $t$ , we plot the best policy performance obtained by the algorithm so far averaged across seeds (note that this quantity is monotonic, since the best policy obtained so far can never be worse at a later point of time when running the algorithm). For BC/BC-O/ILPO, we present a simplified view by picking the best policy obtained through the course of running the algorithm and averaging it across seeds (so the curves are flat lines). As figure 6 shows, MobILE reliably hits expert performance faster than GAIL and GAIFO while often matching/outperforming ILPO/BC/BC-O.

### D.4 Ablation Study on Number of Models used for Strategic Exploration Bonus

In this experiment, we present an ablation study on using more number of models in the ensemble for setting the strategic exploration bonus. Figure 7 suggests that even utilizing two models for purposes of setting the bonus is effective from a practical perspective.

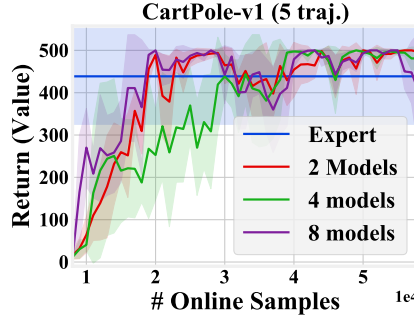


Figure 7: Learning curves for Cartpole-v1 with varying number of dynamics models for assigning bonuses for strategic exploration.