

Active Exploration for Inverse Reinforcement Learning

David Lindner

ETH Zurich

DAVID.LINDNER@INF.ETHZ.CH

Andreas Krause

ETH Zurich

KRAUSEA@ETHZ.CH

Giorgia Ramponi

ETH AI Center

GIORGIA.RAMPONI@AI.ETHZ.CH

Abstract

Inverse Reinforcement Learning (IRL) is a powerful paradigm for inferring a reward function from expert demonstrations. Many IRL algorithms require a known transition model and sometimes even a known expert policy, or they at least require access to a generative model. However, these assumptions are too strong for many real-world applications, where the environment can be accessed only through sequential interaction. We propose a novel IRL algorithm: **Active exploration for Inverse Reinforcement Learning (AceIRL)**, which actively explores an unknown environment and expert policy to quickly learn the expert’s reward function and identify a good policy. AceIRL uses previous observations to construct confidence intervals that capture plausible reward functions and find exploration policies that focus on the most informative regions of the environment. AceIRL is the first approach to active IRL with sample-complexity bounds that does not require a generative model of the environment. AceIRL matches the sample complexity of active IRL with a generative model in the worst case. Additionally, we establish a problem-dependent bound that relates the sample complexity of AceIRL to the suboptimality gap of a given IRL problem. We empirically evaluate AceIRL in simulations and find that it significantly outperforms more naive exploration strategies.

Keywords: Inverse Reinforcement Learning, Active Learning, Reward-free Exploration

1. Introduction

Reinforcement Learning (RL) has achieved impressive results recently, from playing video games [Mnih et al., 2015] to solving robotic control problems [Haarnoja et al., 2019]. However, in many applications, it is challenging to design a reward function that robustly describes the desired task [Amodei et al., 2016]. Instead of using an explicit reward function, Inverse Reinforcement Learning (IRL) seeks to recover the reward by observing an *expert*, e.g., an human who already knows how to perform a task [Ng et al., 2000]. However, most existing IRL algorithms assume that the transition model, and in some cases, the expert’s policy, are *known*. In many real-world applications, this is not given, and the agent needs to estimate the transition dynamics and the expert policy from samples. Figure 1 shows an illustrative example where the agent has to explore an environment and query the expert policy in order to infer the expert’s reward function. IRL with sample-based estimation was only recently analyzed formally by Metelli et al. [2021]. However, they assume a *generative model* of the environment, i.e., the agent can query the transition dynamics for arbitrary

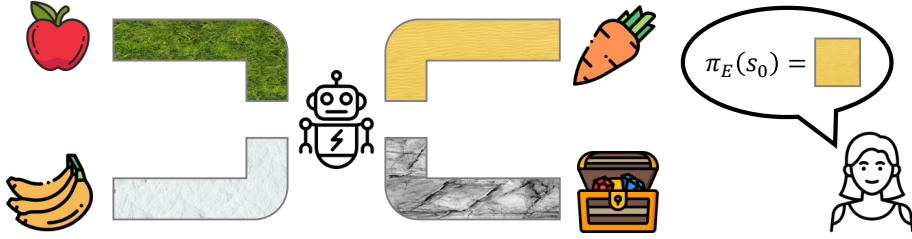


Figure 1: An illustrative example of the Active IRL problem. The agent can choose between four paths that lead to different objects. Observing the expert actions is not enough to infer a reward function (e.g., from observing the expert recommending to take the yellow path, the agent cannot infer that the human prefers to find the carrot). Therefore, the agent has to explore the environment and learn about its dynamics to infer a good reward function.

states and actions. In practice, this assumption is unrealistic, and the agent has to explore the environment from a starting state or state distribution.

In this work, we consider IRL with unknown transition dynamics and expert policy and focus on exploring the environment in order to recover the expert’s reward function efficiently. We propose a novel algorithm, **Active exploration for Inverse Reinforcement Learning (AceIRL)**, which actively explores the environment and the expert policy to infer a good reward function. To the best of our knowledge, we present the first paper providing sample complexity guarantees for the active IRL problem without access to a generative model. The proofs of all results presented in the main paper can be found in Appendix C, we provide more discussion of related work in Appendix B, and in Appendix A we evaluate AceIRL empirically in simulated environments.

2. Preliminaries

Let us first introduce necessary background and notation that we use throughout the paper.

Markov decision process (MDP). A finite-horizon (or episodic) MDP without reward (MDP\R) is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, P, H, s_0)$, where \mathcal{S} is the finite state space of size S ; \mathcal{A} is the finite action space of size A ; $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ is the transition model; H is the horizon and s_0 is the initial state.¹ We describe an agent’s behaviour with a (possible stochastic) policy $\pi \in \mathcal{S} \times [H] \rightarrow \Delta_{\mathcal{A}}$.

Reward function. A reward function $r : \mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, R_{\max}]$ maps state-action-time step triplets to a reward. Given an MDP\R \mathcal{M} and a reward function r , we denote the resulting MDP by $\mathcal{M} \cup r$.

Value functions and optimality conditions. We define the *Q-function* $Q_{\mathcal{M} \cup r}^{\pi, h}(s, a)$ and *value-function* $V_{\mathcal{M} \cup r}^{\pi, h}(s)$ as: $Q_{\mathcal{M} \cup r}^{\pi, h}(s, a) = r_h(s, a) + \sum_{s', a'} \pi_{h+1}(a'|s') P(s'|s, a) Q_{\mathcal{M} \cup r}^{\pi, h+1}(s', a')$, and $V_{\mathcal{M} \cup r}^{\pi, h}(s) = \sum_a \pi_h(a|s) Q_{\mathcal{M} \cup r}^{\pi, h}(s, a)$, respectively.

We define the *advantage function* $A_{\mathcal{M} \cup r}^{\pi, h}(s, a)$ as $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) = Q_{\mathcal{M} \cup r}^{\pi, h}(s, a) - V_{\mathcal{M} \cup r}^{\pi, h}(s)$. A policy π is optimal if $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) \leq 0$ for each time step $h \in [H]$, state $s \in \mathcal{S}$, action $a \in \mathcal{A}$. We denote the set of optimal policies for the MDP $\mathcal{M} \cup r$ with $\Pi_{\mathcal{M} \cup r}^*$.

1. We can model any initial state distribution as a single initial state by modifying the transitions.

3. Active Learning for Inverse Reinforcement Learning (Active IRL)

In this section, we first introduce the Active IRL problem . Then, we define the feasible reward set for finite-horizon MDPs and characterize the error propagation on the reward function and the value function , extending results by Metelli et al. [2021] to finite horizons.

Problem Definition. We want to explore to construct a dataset of demonstrations \mathcal{D} such that an arbitrary IRL algorithm can recover a *good* reward function from it. To be agnostic to the choice of IRL algorithm, we consider the set of all feasible reward functions for a specific expert policy. Formally, we consider IRL problems (\mathcal{M}, π^E) consisting of an MDP \mathcal{M} and an expert policy π^E , and we define the feasible reward set $\mathcal{R}_{\mathcal{M} \cup \pi^E}$ as the set of all reward functions for which π^E is optimal. Let us now define the goal of the active IRL problem formally by providing an optimality criterion.

Definition 1 (Optimality Criterion) Let \mathcal{S} be a sampling strategy. Let $\mathcal{R}_{\mathfrak{B}}$ be the exact feasible set and $\mathcal{R}_{\hat{\mathfrak{B}}}$ be the feasible set recovered after observing $n \geq 0$ samples collected from \mathcal{M} and π^E . We say that \mathcal{S} is (ϵ, δ, n) -correct if with probability at least $1 - \delta$ it holds that:

$$\inf_{\hat{r} \in \mathcal{R}_{\hat{\mathfrak{B}}}} \sup_{\hat{\pi}^* \in \Pi_{\widehat{\mathcal{M}} \cup \hat{r}}^*} \max_{s,a,h} \left| Q_{\mathcal{M} \cup r}^{\pi^*,h}(s,a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*,h}(s,a) \right| \leq \epsilon \quad \text{for each } r \in \mathcal{R}_{\mathfrak{B}},$$

$$\inf_{r \in \mathcal{R}_{\mathfrak{B}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup r}^*} \max_{s,a,h} \left| Q_{\mathcal{M} \cup r}^{\pi^*,h}(s,a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*,h}(s,a) \right| \leq \epsilon \quad \text{for each } \hat{r} \in \mathcal{R}_{\hat{\mathfrak{B}}},$$

where π^* is an optimal policy in $\mathcal{M} \cup r$ and $\hat{\pi}^*$ is an optimal policy in $\widehat{\mathcal{M}} \cup \hat{r}$.

The first condition states that for each reward in the exact feasible set, the best reward we could estimate in the recovered feasible set has a low error everywhere. This condition is a type of “recall”: every possible true reward function needs to be captured by the recovered feasible set. The second condition ensures that there is a possible true reward function with a low error for every possible recovered reward function. This avoids an unnecessarily large recovered feasible set. This condition is a type of “precision”: if we recover a reward function, it has to be close to a possible true reward function.

Feasible Rewards in Finite-horizon MDPs Ng et al. [2000] characterize the feasible reward set implicitly in the infinite horizon setting, whereas Metelli et al. [2021] characterize it explicitly. Here, we provide similar results for a finite horizon.

Lemma 2 (Feasible Reward Set Implicit) A reward function r is feasible if and only if for all s, a, h it holds that: $A_{\mathcal{M} \cup r}^{\pi,h}(s,a) = 0$ if $\pi_h^E(a|s) \geq 0$ and $A_{\mathcal{M} \cup r}^{\pi,h}(s,a) \leq 0$ if $\pi_h^E(a|s) = 0$. Moreover, if the second inequality is strict, π^E is uniquely optimal, i.e., $\Pi_{\mathcal{M} \cup r}^* = \{\pi^E\}$.

Lemma 3 (Feasible Reward Set Explicit) A reward function r is feasible if and only if there exists an $\{A_h \in \mathbb{R}_{\geq 0}^{S \times \mathcal{A}}\}_{h \in [H]}$ and $\{V_h \in \mathbb{R}^S\}_{h \in [H]}$ such that for all s, a, h it holds that:

$$r_h(s,a) = -A_h(s,a) \mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s) + \sum_{s'} P(s'|s,a) V_{h+1}(s')$$

Here, the **first term** ensures there is an advantage function for π^E and it is 0 for actions the expert takes and $A_h(s,a)$ for actions the expert does not take. The **second term** corresponds to reward-shaping by the value function.

Algorithm 1 AceIRL algorithm for IRL in an unknown environment.

- 1: **Input:** significance $\delta \in (0, 1)$, accuracy ϵ , IRL algorithm \mathcal{A} , number of episodes N_E
 - 2: Initialize $k \leftarrow 0$, $\epsilon_0 \leftarrow H/10$
 - 3: **while** $\epsilon_k > \epsilon/4$ **do**
 - 4: Solve (convex) optimization problem (ACE) to obtain π_k
 - 5: Explore with policy π_k for N_E episodes, observing transitions and expert actions
 - 6: Increment $k \leftarrow k + 1$ and update $\hat{P}_k, \hat{\pi}_k, C_k^h$, and $\hat{r}_k \leftarrow \mathcal{A}(\mathcal{R}_{\mathfrak{B}})$
 - 7: Update accuracy $\epsilon_k \leftarrow \max_a \hat{E}_k^0(s_0, a)$
 - 8: **end while**
 - 9: **return** Estimated reward function \hat{r}_k
-

Error Propagation Next, we study the error propagation of estimating the transition model P with \hat{P} and the expert policy π^E with $\hat{\pi}^E$. In particular, we bound the estimation error on the reward as a function of the estimation errors of \hat{P} and $\hat{\pi}^E$, extending a result by Metelli et al. [2021] to the finite horizon setting.

Theorem 4 (Error Propagation) *Let (\mathcal{M}, π^E) and $(\widehat{\mathcal{M}}, \widehat{\pi}^E)$ be two IRL problems. Then, for any $r \in \mathcal{R}_{(\mathcal{M}, \pi^E)}$ there exists $\hat{r} \in \widehat{\mathcal{R}}_{(\widehat{\mathcal{M}}, \widehat{\pi}^E)}$ such that:*

$$|r_h(s, a) - \hat{r}_h(s, a)| \leq A_h(s, a) |\pi_h^E(a|s) - \hat{\pi}_h^E(a|s)| + \sum_{s'} V_{h+1}(s') |P(s'|s, a) - \hat{P}(s'|s, a)|$$

and we can bound $V_h \leq (H - h)R_{\max}$ and $A_h \leq (H - h)R_{\max}$.

4. Active Exploration for Inverse Reinforcement Learning

Let us now turn to our original problem: recovering the expert’s reward function in an unknown environment *without* a generative model. This problem is harder since we need to create an exploration strategy to acquire the desired information about the environment. To address this problem, we propose **Active** exploration for **Inverse Reinforcement Learning** (AceIRL). First, we introduce a simplified version of the algorithm, which comes with a problem independent sample complexity result (Section 4.1). Next, we introduce the full algorithm, which considers the expected reduction of uncertainty in the next iteration to improve exploration and maintains a confidence set of plausibly optimal policies to focus on the most relevant regions (Section 4.2). The full algorithm provides a tighter, problem-dependent sample complexity bound (Section 4.3).

4.1 Uncertainty-based Exploration for IRL

The first idea of AceIRL is similar to reward-free UCRL [Kaufmann et al., 2021]. Our goal is to fulfill the PAC requirement in Definition 1. Hence, we start from an upper bound on the estimation error between the performance of the optimal policy $\hat{\pi}^*$ for a reward $\hat{r} \in \mathcal{R}_{\mathfrak{B}}$ in the recovered feasible set and the optimal policy π^* for a reward function $r \in \mathcal{R}_{\mathfrak{B}}$ in the true MDP \mathcal{M} . We will then use this upper bound to drive the exploration. For each timestep h and iteration k , we define the error: $\hat{e}_k^h(s, a; \pi^*, \hat{\pi}^*) = \left| Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) \right|$.

We can define an upper bound on these errors recursively with $C_k^H(s, a) = 0$ and

$$E_k^h(s, a) = \min\left((H - h)R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a' \in \mathcal{A}} E_k^{h+1}(s', a')\right). \quad (\text{EB1})$$

It is straightforward to show that $\hat{e}_k^h(s, a; \pi^*, \hat{\pi}^*) \leq E_k^h(s, a)$ for any two policies $\pi^*, \hat{\pi}^*$. Using this error bound, we can introduce a simplified version of AceIRL that explores greedily with respect to $E_k^h(s, a)$. We call this algorithm ‘‘AceIRL Greedy’’. Note that this is equivalent to solving the RL problem defined by $\mathcal{M} \cup C_k^h$; hence, we can use any RL solver to find the exploration policy in practice. We can show that if we stop if $4 \max_a E_k^0(s_0, a) \leq \epsilon$, the solution fulfills the PAC requirement in Definition 1. Furthermore, we show in Appendix C.4 that AceIRL Greedy achieves a sample complexity on order $\tilde{O}(H^5 R_{\max}^2 SA/\epsilon^2)$, which matches the sample complexity of uniform sampling *with a generative model*. This result implies that we do not need a generative model to achieve a good sample complexity for IRL.

4.2 Problem Dependent Exploration

AceIRL Greedy is limited in two ways: (i) it explores states that have high uncertainty so far, whereas our goal is to reduce uncertainty *in the next iteration*, and (ii) it explores to reduce the uncertainty about all policies, whereas our goal is to reduce the uncertainty primarily about *plausibly optimal* policies. To address these limitations, we propose two modifications that yield the full AceIRL algorithm.

Reducing future uncertainty. The greedy policy w.r.t. E_k^h explores states in which the estimation error on the Q-functions is large. But this is not exactly what reduces our uncertainty the most. Ideally, we would explore with a policy that minimizes E_{k+1}^h . However, we cannot compute this quantity exactly. Instead, we can approximate it using our current estimate of the transition model. Concretely, if we have an exploration policy π , we can estimate the reward uncertainty at the next iteration as $\hat{C}_{k+1}^h(s, a) = (H - h)R_{\max} \min\left(1, 2\sqrt{\frac{2\ell_k^h(s, a)}{n_k^h(s, a) + \hat{n}_\pi^h(s, a)}}\right)$, where $\hat{n}_\pi^h(s, a) = N_E \cdot \eta_{\mathcal{M}, \pi}^{0, h}(s, a|s_0)$ is the expected number of times π visits (s, a) at time h and N_E is the number of episodes we will explore with π . We can use this estimate to find a policy that minimizes our estimate of E_{k+1}^h . While our original approach was akin to ‘‘uncertainty sampling’’, we now have a better way to measure the ‘‘informativeness’’ of choosing an exploration policy.

Focusing on plausibly optimal policies. By exploring greedily w.r.t. E_k^h , we reduce the estimation error of all policies. However, we are primarily interested in estimating the distance between policies $\pi^* \in \Pi_{\mathcal{M} \cup r}^*$ and $\hat{\pi}^* \in \Pi_{\mathcal{M} \cup \hat{r}}^*$ with $r \in \mathcal{R}_{\mathfrak{B}}$ and $\hat{r} \in \mathcal{R}_{\hat{\mathfrak{B}}}$. Of course, we do not know these sets. Instead, assume we can construct a set of plausibly optimal policies $\hat{\Pi}_k$ that contains all π^* and $\hat{\pi}^*$ with high probability. Then, we can redefine our upper bounds recursively as $\hat{E}_k^H(s, a) = 0$ and:

$$\hat{E}_k^h(s, a) = \min\left((H - h)R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{\pi \in \hat{\Pi}_{k-1}} \pi(a'|s') \hat{E}_k^{h+1}(s', a')\right), \quad (\text{EB2})$$

In contrast to (EB1), we maximize over policies in $\hat{\Pi}_k$ rather than all actions. Acting greedily with respect to $\hat{E}_k^h(s, a)$ is equivalent to finding the optimal policy $\pi_k \in \hat{\Pi}_k$ for

the RL problem defined by $\mathcal{M} \cup C_k^h$. To construct the set of plausibly optimal policies, we use an arbitrary IRL algorithm \mathcal{A} . We only assume that \mathcal{A} will return a reward function $\hat{r}_k \in \mathcal{R}_{\mathfrak{B}}$. Then, we can construct a set of plausibly optimal policies as $\hat{\Pi}_k = \{\pi | V_{\mathcal{M} \cup \hat{r}_k}^{\pi^*}(s_0) - V_{\mathcal{M} \cup \hat{r}_k}^{\pi}(s_0) \leq 10\epsilon_k\}$. We show in Appendix C.5 that $\hat{\Pi}_k$ contains both π^* and $\hat{\pi}_k^*$ with high probability. We can define a stopping condition analogously to before: $4 \max_a \hat{E}_k^0(s_0, a) \leq \epsilon$. Again, we can prove that if the algorithm stops, then $\mathcal{R}_{\mathfrak{B}}$ respects Definition 1.

Implementing AceIRL. To implement the full algorithm, we need to solve:

$$\pi_k \in \underset{\pi}{\operatorname{argmin}} \max_{\hat{\pi} \in \hat{\Pi}_{k-1}} \hat{E}_{k+1}^0(s_0, \hat{\pi}(s_0)) \quad (\text{ACE})$$

The policy solving this problem minimizes the uncertainty at the next iteration about plausibly optimal policies. This combines both modifications we just discussed. This problem might seem difficult to solve at first, but, perhaps surprisingly, it can be formulated as a convex optimization problem solvable with standard techniques (cf. Appendix C.6).

4.3 Sample Complexity of AceIRL

In this section, we present our main result about the sample complexity of AceIRL. The result is problem-dependent and depends on the advantage function $A_{\mathcal{M} \cup r}^{*,h}(s, a)$, which acts similarly to a suboptimality gap: the closer the value of the second best action is to the best action, the harder it is to identify the best action and infer the correct reward function.

Theorem 5 [AceIRL Sample Complexity] *AceIRL returns a (ϵ, δ, n) -correct solution with*

$$n \leq \tilde{\mathcal{O}} \left(\min \left[\frac{H^5 R_{\max}^2 S A}{\epsilon^2}, \frac{H^4 R_{\max}^2 S A \epsilon_{\tau-1}^2}{\min_{s,a,h} (A_{\mathcal{M} \cup r}^{*,h}(s, a))^2 \epsilon^2} \right] \right)$$

where $\epsilon_{\tau-1}$ depends on the choice of N_E , the number of episodes of exploration in each iteration. $A_{\mathcal{M} \cup r}^{*,h}(s, a)$ is the advantage function of $r \in \underset{r \in \mathcal{R}_{\mathfrak{B}}}{\operatorname{argmin}} \max_{h,s,a} (r_h(s, a) - \hat{r}_{k,h}(s, a))$, the reward function from the feasible set $\mathcal{R}_{\mathfrak{B}}$ closest to the estimated reward function \hat{r}_k .

This result is the minimum of two terms. The first term is problem independent and it is achieved both by AceIRL Greedy and the full AceIRL. Using (ACE) can yield a better sample complexity, represented by the second term in the minimum. This bound depends on two main components: the ratio $\epsilon_{\tau-1}/\epsilon$ and the advantage function $A_{\mathcal{M} \cup r}^{*,h}(s, a)$. The ratio depends on the choice of N_E , the number of exploration episodes per iteration. Appendix C.5 provides the full proof of this theorem.

5. Conclusion

We considered active inverse reinforcement learning (IRL) with unknown transition dynamics and expert policy and introduced AceIRL, an efficient exploration strategy to learn about both the dynamic and the expert policy. We make a crucial step towards IRL algorithms with theoretical guarantees, but there remain many questions for future work, such as extending the approach to continuous environments using function approximation.

References

- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. *arXiv:1606.06565*, 2016.
- Daniel S Brown, Yuchen Cui, and Scott Niekum. Risk-aware active inverse reinforcement learning. In *Conference on Robot Learning (CoRL)*, 2018.
- Robert Cohn, Edmund Durfee, and Satinder Singh. Comparing action-query strategies in semi-autonomous agents. In *AAAI Conference on Artificial Intelligence*, 2011.
- Gregory Dexter, Kevin Bello, and Jean Honorio. Inverse reinforcement learning in a continuous state space with formal guarantees. In *Advances in Neural Information Processing Systems*, 2021.
- Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning to walk via deep reinforcement learning. In *Proceedings of Robotics: Science and Systems (RSS)*, 2019.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. In *Proceedings of International Conference on Machine Learning (ICML)*, 2020.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of International Conference on Machine Learning (ICML)*, 2002.
- Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Edouard Leurent, and Michal Valko. Adaptive reward-free exploration. In *Algorithmic Learning Theory*, 2021.
- Johannes Kulick, Marc Toussaint, Tobias Lang, and Manuel Lopes. Active learning for teaching a robot grounded relational symbols. In *Proceedings of International Joint Conferences on Artificial Intelligence*, 2013.
- Sergey Levine, Zoran Popovic, and Vladlen Koltun. Nonlinear inverse reinforcement learning with Gaussian processes. *Advances in Neural Information Processing Systems*, 2011.
- Manuel Lopes, Francisco Melo, and Luis Montesano. Active learning for reward estimation in inverse reinforcement learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2009.
- Dylan P Losey and Marcia K O’Malley. Including uncertainty when learning from human corrections. In *Conference on Robot Learning (CoRL)*, 2018.

- Alberto Maria Metelli, Giorgia Ramponi, Alessandro Concetti, and Marcello Restelli. Provably efficient learning of transferable rewards. In *Proceedings of International Conference on Machine Learning (ICML)*, 2021.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *Proceedings of International Conference on Machine Learning (ICML)*, 2000.
- Brendan O’Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd. Conic optimization via operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and Applications*, 169(3):1042–1068, June 2016. URL <http://stanford.edu/~boyd/papers/scs.html>.
- Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *Proceedings of International Joint Conferences on Artificial Intelligence*, 2007.
- Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. Maximum margin planning. In *Proceedings of International Conference on Machine Learning (ICML)*, 2006.
- Andrea Zanette, Mykel J Kochenderfer, and Emma Brunskill. Almost horizon-free structure-aware best policy identification with a generative model. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, 2019.
- Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2008.

Appendix

Table of Contents

A	Simulation Experiments	9
B	Related Work	11
C	Proofs of Theoretical Results	11
C.1	Simulation Lemmas	11
C.2	Feasible Reward Set	16
C.3	Uniform Sampling IRL with a Generative Model	18
C.4	Sample Complexity of AceIRL in Unknown Environments (Problem Independent)	22
C.5	Sample Complexity of AceIRL in Unknown Environments (Problem Dependent)	28
C.6	Computing the Exploration Policy	34
D	Experimental Details	35
D.1	Details on the Environments	36
D.2	Implementation Details	36
D.3	Additional Results	37
E	Connection to Reward-free Exploration	38

Appendix A. Simulation Experiments

We perform a series of simulation experiments to evaluate AceIRL. We simulate a (deterministic) expert policy using an underlying true reward function, and compare it to the recovered reward functions. Our main evaluation metric is a *normalized regret*:

$$(V_{\mathcal{M} \cup r}^{\pi^*, 0}(s_0) - V_{\mathcal{M} \cup r}^{\hat{\pi}^*, 0}(s_0)) / (V_{\mathcal{M} \cup r}^{\pi^*, 0}(s_0) - V_{\mathcal{M} \cup r}^{\bar{\pi}^*, 0}(s_0)),$$

where π^* is the optimal policy for $\mathcal{M} \cup r$, $\hat{\pi}^*$ is the optimal policy for $\widehat{\mathcal{M}} \cup \hat{r}$, and $\bar{\pi}^*$ is the worst possible policy for r , i.e., the optimal policy for $\mathcal{M} \cup (-r)$.

We introduce the *Four Paths* environment shown in Figure 1, which consists of four chains of states that have different randomly sampled transition probabilities. One path has a goal with reward 1; all other rewards are 0. We also evaluate on *Random MDPs* with uniformly sampled transition dynamics and reward functions, the *Double Chain* environment proposed by Kaufmann et al. [2021], and the *Chain* and *Gridworld* environments proposed by Metelli et al. [2021]. Appendix D.1 provides details on the transition dynamics and rewards of all environments.

	Uniform sampling (gener. model)	TRAVEL (gener. model) [Metelli et al., 2021]	Random Exploration	AceIRL Greedy	AceIRL (Full)
Four Paths (Figure 1)	1900 ± 71		17840 ± 1886		
– $N_E = 50$		1560 ± 76		24180 ± 1747	10780 ± 1369
– $N_E = 100$		2000 ± 0		32760 ± 2172	14080 ± 1603
– $N_E = 200$		4000 ± 0		52000 ± 4057	16160 ± 2033
Double Chain [Kaufmann et al., 2021]	1980 ± 66		23640 ± 2195		
– $N_E = 50$		1120 ± 46		16240 ± 842	11580 ± 870
– $N_E = 100$		2000 ± 0		22200 ± 1329	15440 ± 1031
– $N_E = 200$		4000 ± 0		37200 ± 1664	20400 ± 1629
Metelli et al. [2021]:					
Random MDPs ($N_E = 1$)	22 ± 1	27 ± 1	22 ± 1	23 ± 1	21 ± 1
Chain ($N_E = 1$)	78 ± 2	76 ± 4	161 ± 8	153 ± 8	142 ± 9
Gridworld ($N_E = 1$)	43 ± 2	35 ± 2	45 ± 2	46 ± 3	48 ± 2

Table A.1: Sample complexity of AceIRL compared to random exploration and methods that use a generative model. We show the number of samples necessary to find a policy with normalized regret less than 0.4. We report means and standard errors computed over 50 random seeds each. For each environment, we highlight in **bold** the method that achieves the best performance without access to a generative model. If multiple methods are within one standard error distance, we highlight all of them. Overall, AceIRL is the most sample efficient method without a generative model if N_E is chosen small enough. In Appendix D.3, we show learning curves for all individual experiments.

We compare AceIRL and AceIRL Greedy to a uniformly random exploration policy, as a naive exploration strategy. Further, we consider uniform sampling with a generative model as well as TRAVEL [Metelli et al., 2021], which can be more sample efficient because they do not need to explore the environment. Note that TRAVEL is designed to learn a reward to be transferred to a known target environment. Instead, we use a modified version that uses the estimated MDP as a target. Appendix D.2 provides more details on our implementations, and we provide source code in the supplemental material.

Table A.1 shows the sample efficiency of all algorithms in all environments, measured as the number of samples needed to achieve a regret threshold of 0.4 (different thresholds yield similar conclusions; cf. Appendix D). AceIRL is the most sample efficient exploration strategy without access to a generative model; but the relative

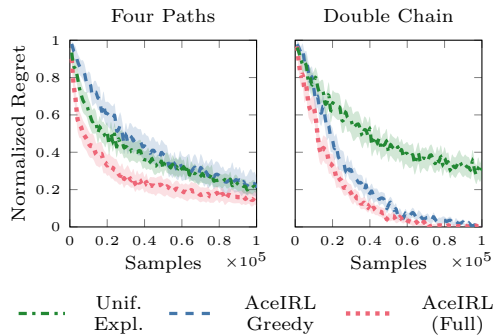


Figure A.1: Normalized regret (lower is better) of the policy optimizing for the inferred reward in the estimated MDP as a function of the number of samples. The plots show the mean and 95% confidence intervals computed using 50 random seeds. We use $N_E = 50$.

differences between the methods depend on the environment. In some cases, AceIRL even performs comparably to methods using a generative model, as the theory predicts.

In the *Four Paths* and *Double Chain* environments, we also vary the N_E parameter. AceIRL performs better for small values at the computational cost of updating the exploration policy more often. If N_E is too large, using AceIRL can be as bad as a uniformly random exploration policy. Increasing N_E hurts the performance of AceIRL Greedy more severely, which does not consider N_E explicitly. Figure A.1 shows the normalized regret as a function of the number of samples in *Four Paths* and *Double Chain*. In both cases AceIRL performs best. However, AceIRL Greedy is worse than random exploration in the *Four Paths* environment. Hence, we find that the problem dependent exploration strategy of the full algorithm significantly improves the sample efficiency.

Appendix B. Related Work

Most IRL algorithms assume that the underlying transition model is known [Ratliff et al., 2006, Ziebart et al., 2008, Ramachandran and Amir, 2007, Levine et al., 2011]. However, the transition model usually needs to be estimated from samples, which induces an error in the recovered reward function that most papers do not study. Metelli et al. [2021] analyze this error and the sample complexity of IRL in a tabular setting with a generative model. They propose an algorithm focused on transferring the learned reward function to a fully known target environment. Dexter et al. [2021] provides a similar analysis in continuous state spaces and discrete action spaces, but they still require a generative model of the environment. In contrast, we *do not* assume access to a generative model and thus need to tackle the exploration problem in IRL.

Some prior work studies active learning algorithms for IRL in a Bayesian framework but without theoretical guarantees. Lopes et al. [2009] propose an active learning algorithm for IRL that estimates a posterior distribution over reward functions from demonstrations, requiring a prior distribution and full knowledge of the environment dynamics. Relatedly, Cohn et al. [2011] consider a Bayesian IRL setting with a semi-autonomous agent that asks an expert for advice if it is uncertain about the reward. Brown et al. [2018] empirically study active IRL in several safety-critical environments, selecting queries using value at risk. Kulick et al. [2013] consider active learning for a robotic manipulation task, asking a human expert for advice in situations with the highest predictive uncertainty. Similarly, Losey and O’Malley [2018] propose a method to learn uncertainty estimates from human corrections in a robotics context. All of these papers assume a Bayesian framework and do not provide theoretical guarantees. In contrast, our setup does not require a prior over reward functions, and we provide theoretical sample complexity guarantees for our algorithm.

Appendix C. Proofs of Theoretical Results

In Table C.2, we provide a reference of the notation and symbols used in our paper.

C.1 Simulation Lemmas

In this section, we establish several simulation lemmas that we will use throughout our analysis. Some of the results were already derived in prior work for the infinite horizon

Table C.2: Overview of our notation

Symbol	Name	Signature
\mathcal{M}	Markov decision process without reward (MDP\R)	$(\mathcal{S}, \mathcal{A}, P, H, s_0)$
\mathcal{S}	State space	
\mathcal{A}	Action space	
P	Transition model	$\mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$
H	Horizon	$H \in \mathbb{N}^+$
s_0	Initial state	$s_0 \in \mathcal{S}$
π	Policy	$\mathcal{S} \times [H] \rightarrow \Delta_{\mathcal{A}}$
r	Reward function	$\mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, R_{\max}], R_{\max} \in \mathbb{R}^+$
$\mathcal{M} \cup r$	Markov decision process (MDP)	$(\mathcal{S}, \mathcal{A}, P, H, s_0, r)$
$Q_{\mathcal{M} \cup r}^{\pi, h}$	Q-function (of π in $\mathcal{M} \cup r$)	$\mathcal{S} \times \mathcal{A} \times [H] \rightarrow \mathbb{R}$
$V_{\mathcal{M} \cup r}^{\pi, h}$	Value function (of π in $\mathcal{M} \cup r$)	$\mathcal{S} \times [H] \rightarrow \mathbb{R}$
$A_{\mathcal{M} \cup r}^{\pi, h}$	Advantage function (of π in $\mathcal{M} \cup r$)	$\mathcal{S} \times \mathcal{A} \times [H] \rightarrow \mathbb{R}$
$\eta_{\mathcal{M}, \pi}^{h, \cdot}(\cdot s_0)$	State-visitation frequency (conditioned on state)	$[H] \rightarrow \Delta_{\mathcal{S}}$
$\eta_{\mathcal{M}, \pi}^{h, \cdot}(\cdot s_0, a_0)$	State-visitation frequency (conditioned on state-action)	$[H] \rightarrow \Delta_{\mathcal{S}}$
$\eta_{\mathcal{M}, \pi}^{h, \cdot}(\cdot, \cdot s_0)$	State-action-visitation frequency (conditioned on state)	$[H] \times \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$
$\eta_{\mathcal{M}, \pi}^{h, \cdot}(\cdot, \cdot s_0, a_0)$	State-action-visitation frequency (conditioned on state)	$[H] \times \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$
$\mathcal{R}_{\mathcal{M} \cup r}$	Feasible set of $\mathcal{M} \cup r$	
$\mathcal{R}_{\mathfrak{B}} = \mathcal{R}_{\mathcal{M} \cup \pi^E}$	Exact feasible set	
$\mathcal{R}_{\hat{\mathfrak{B}}} = \mathcal{R}_{\widehat{\mathcal{M} \cup \hat{\pi}^E}}$	Recovered feasible set	
ϵ	Target accuracy	$\epsilon \in \mathbb{R}^+$
δ	Significancy	$\delta \in (0, 1)$
N_E	Number of exploration episodes	$N_E \in \mathbb{N}^+$

setting, e.g., by Zanette et al. [2019] and Metelli et al. [2021]. For completeness, we provide proofs for all results in the finite-horizon setting.

Definition C.1 (Occupancy measures) We define $\eta_{\mathcal{M}, \pi}^{h, h'}(s | s_0)$ as the probability of being in state s at timestep $h' \geq h$ following a policy π in MDP\R \mathcal{M} starting in state s_0 at timestep h . We can compute it recursively as:

$$\begin{aligned} \eta_{\mathcal{M}, \pi}^{h, h}(s' | s) &:= \mathbb{1}_{\{s'=s\}} \\ \eta_{\mathcal{M}, \pi}^{h, h'+1}(s' | s) &:= \sum_{s'', \tilde{a}} P(s' | s'', \tilde{a}) \pi_{h'}(\tilde{a} | s'') \eta_{\mathcal{M}, \pi}^{h, h'}(s'' | s) \end{aligned}$$

We define the same probability for state-action pairs analogously:

$$\begin{aligned} \eta_{\mathcal{M}, \pi}^{h, h'}(s', a' | s, a) &:= \mathbb{1}_{\{s'=s, a'=a\}} \\ \eta_{\mathcal{M}, \pi}^{h, h'+1}(s', a' | s, a) &:= \sum_{\tilde{s}, \tilde{a}} \pi_{h'}(a' | s') P(s' | \tilde{s}, \tilde{a}) \eta_{\mathcal{M}, \pi}^{h, h'}(\tilde{s}, \tilde{a} | s, a) \end{aligned}$$

as well as

$$\begin{aligned} \eta_{\mathcal{M}, \pi}^{h, h}(s', a' | s) &:= \pi_h(a' | s') \mathbb{1}_{\{s'=s\}} \\ \eta_{\mathcal{M}, \pi}^{h, h'+1}(s', a' | s) &:= \sum_{\tilde{s}, \tilde{a}} \pi_{h'}(a' | s') P(s' | \tilde{s}, \tilde{a}) \eta_{\mathcal{M}, \pi}^{h, h'}(\tilde{s}, \tilde{a} | s) \end{aligned}$$

Because the environment is Markovian, it also holds for $h' > h$ that

$$\eta_{\mathcal{M},\pi}^{h,h'}(s'|s) = \sum_{\tilde{s},a} \eta_{\mathcal{M},\pi}^{h+1,h'}(s'|\tilde{s})P(\tilde{s}|s,a)\pi_h(a|s)$$

and equivalently for state-action pairs.

Lemma C.2 *The value function and Q-function of a policy π in an MDP $\mathcal{M}_{\cup r}$ at timestep h can be expressed as:*

$$V_{\mathcal{M}_{\cup r}}^{\pi,h}(s) = \sum_{h'=h}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{h,h'}(s',a'|s)r_{h'}(s',a')$$

$$Q_{\mathcal{M}_{\cup r}}^{\pi,h}(s,a) = \sum_{h'=h}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{h,h'}(s',a'|s,a)r_{h'}(s',a')$$

Proof We show the result for the value function; the derivation for the Q-function is analogous.

Note that for $h = H$ the statement holds because $V_{\mathcal{M}_{\cup r}}^{\pi,H}(s) = 0$. The general result follows by induction. Assume that for $h + 1$ the statement holds. Then:

$$\begin{aligned} V_{\mathcal{M}_{\cup r}}^{\pi,h}(s) &\stackrel{(a)}{=} \sum_a \pi_h(a|s) \left(r_h(s,a) + \sum_{s'} P(s'|s,a) V_{\mathcal{M}_{\cup r}}^{\pi,h+1}(s') \right) \\ &\stackrel{(b)}{=} \sum_a \pi_h(a|s) \left(r_h(s,a) + \sum_{s'} P(s'|s,a) \left(\sum_{h'=h+1}^H \sum_{s'',a''} \eta_{\mathcal{M},\pi}^{h+1,h'}(s'',a''|s') r_{h'}(s'',a'') \right) \right) \\ &\stackrel{(c)}{=} \sum_a \pi_h(a|s) r_h(s,a) + \sum_{h'=h+1}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{h,h'}(s'|s) \pi_{h'}(a'|s') r_{h'}(s',a') \\ &\stackrel{(d)}{=} \sum_{h'=h}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{h,h'}(s'|s) \pi_{h'}(a'|s') r_{h'}(s',a') \end{aligned}$$

where (a) uses the Bellman equation, (b) the induction step, (c) uses Theorem C.1 and relabels $s'' \rightarrow s'$, $a'' \rightarrow a'$, and (d) uses Theorem C.1 again and relabels $a \rightarrow a'$. \blacksquare

Lemma C.3 (Simulation lemma 1 by Metelli et al. [2021]) *Let \mathcal{M} be an MDP $\setminus R$, and r, \hat{r} two reward functions with corresponding optimal policies $\pi^*, \hat{\pi}^*$. Then,*

$$Q_{\mathcal{M}_{\cup r}}^{\pi^*,h}(s,a) - Q_{\mathcal{M}_{\cup \hat{r}}}^{\hat{\pi}^*,h}(s,a) \leq \sum_{h'=h}^H \sum_{s',a'} \eta_{\mathcal{M},\pi^*}^{h,h'}(s',a'|s,a) (r_{h'}(s',a') - \hat{r}_{h'}(s',a'))$$

$$V_{\mathcal{M}_{\cup r}}^{\pi^*,h}(s) - V_{\mathcal{M}_{\cup \hat{r}}}^{\hat{\pi}^*,h}(s) \leq \sum_{h'=h}^H \sum_{s',a'} \eta_{\mathcal{M},\pi^*}^{h,h'}(s',a'|s) (r_{h'}(s',a') - \hat{r}_{h'}(s',a'))$$

Proof Note that $Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) \geq Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a)$ for all s, a because $\hat{\pi}^*$ is optimal for \hat{r} . Hence

$$\begin{aligned} Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) &\leq Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\pi^*, h}(s, a) \\ &\stackrel{(a)}{=} \sum_{h'=h}^H \sum_{s', a'} \eta_{\mathcal{M}, \pi^*}^{h, h'}(s', a' | s, a) (r_{h'}(s', a') - \hat{r}_{h'}(s', a')) \end{aligned}$$

where (a) uses Theorem C.2. After observing $V_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s) \geq V_{\mathcal{M} \cup r}^{\pi^*, h}(s)$, the second result follows analogously. \blacksquare

Lemma C.4 *Let \mathcal{M} be an MDP $\setminus R$, r, \hat{r} two reward functions with optimal policies $\pi^*, \hat{\pi}^*$. Then,*

$$Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a) \leq \sum_{h'=h}^H \sum_{s', a'} \left(\eta_{\mathcal{M}, \pi^*}^{h, h'}(s', a' | s, a) - \eta_{\mathcal{M}, \hat{\pi}^*}^{h, h'}(s', a' | s, a) \right) (r_{h'}(s', a') - \hat{r}_{h'}(s', a'))$$

Proof

$$\begin{aligned} Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a) &= (Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a)) + (Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)) \\ &\stackrel{(a)}{\leq} \sum_{h'=h}^H \sum_{s', a'} \eta_{\mathcal{M}, \pi^*}^{h, h'}(s', a' | s, a) (r_{h'}(s', a') - \hat{r}_{h'}(s', a')) + (Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)) \\ &\stackrel{(b)}{=} \sum_{h'=h}^H \sum_{s', a'} \eta_{\mathcal{M}, \pi^*}^{h, h'}(s', a' | s, a) (r_{h'}(s', a') - \hat{r}_{h'}(s', a')) + \sum_{h'=h}^H \sum_{s', a'} \eta_{\mathcal{M}, \hat{\pi}^*}^{h, h'}(s', a' | s, a) (\hat{r}_{h'}(s', a') - r_{h'}(s', a')) \\ &= \sum_{h'=h}^H \sum_{s', a'} \left(\eta_{\mathcal{M}, \pi^*}^{h, h'}(s', a' | s, a) - \eta_{\mathcal{M}, \hat{\pi}^*}^{h, h'}(s', a' | s, a) \right) (r_{h'}(s', a') - \hat{r}_{h'}(s', a')) \end{aligned}$$

where (a) uses Theorem C.3 and (b) uses Theorem C.2. \blacksquare

Lemma C.5 *Let $\mathcal{M}_1, \mathcal{M}_2$ be two MDP $\setminus R$ with transition dynamics P_1, P_2 respectively, r a reward function and π a policy. Then, for any state s and timestep h :*

$$\begin{aligned} V_{\mathcal{M}_2 \cup r}^{\pi, h}(s) - V_{\mathcal{M}_1 \cup r}^{\pi, h}(s) &= \sum_{h'=h}^H \sum_{s', a', s''} \eta_{\mathcal{M}_2, \pi}^{h, h'}(s'; s) \pi_{h'}(a' | s') (P_2(s'' | s', a') - P_1(s'' | s', a')) V_{\mathcal{M}_1 \cup r}^{\pi, h'+1}(s'') \\ V_{\mathcal{M}_1 \cup r}^{\pi, h}(s) - V_{\mathcal{M}_2 \cup r}^{\pi, h}(s) &= \sum_{h'=h}^H \sum_{s', a', s''} \eta_{\mathcal{M}_2, \pi}^{h, h'}(s'; s) \pi_{h'}(a' | s') (P_1(s'' | s', a') - P_2(s'' | s', a')) V_{\mathcal{M}_1 \cup r}^{\pi, h'+1}(s'') \end{aligned}$$

Moreover,

$$\left| V_{\mathcal{M}_2 \cup r}^{\pi, h}(s) - V_{\mathcal{M}_1 \cup r}^{\pi, h}(s) \right| \leq \sum_{h'=h}^H \sum_{s', a', s''} \eta_{\mathcal{M}_2, \pi}^{h, h'}(s'; s) \pi_{h'}(a' | s') \left| P_2(s'' | s', a') - P_1(s'' | s', a') \right| V_{\mathcal{M}_1 \cup r}^{\pi, h'+1}(s'')$$

Proof We start by writing explicitly the value-functions:

$$\begin{aligned} V_{\mathcal{M}_2 \cup r}^{\pi, h}(s) - V_{\mathcal{M}_1 \cup r}^{\pi, h}(s) &= \sum_{a, s'} \pi_h(a|s) \left(P_2(s'|s, a) V_{\mathcal{M}_2 \cup r}^{\pi, h+1}(s') - P_1(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s') \pm P_2(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s') \right) \\ &= \sum_{a, s'} \pi_h(a|s) \left((P_2(s'|s, a) - P_1(s'|s, a)) V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s') + P_2(s'|s, a) (V_{\mathcal{M}_2 \cup r}^{\pi, h+1}(s') - V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s')) \right) \end{aligned}$$

Unrolling the recursion gives the first result; the second result follows similarly:

$$\begin{aligned} V_{\mathcal{M}_1 \cup r}^{\pi, h}(s) - V_{\mathcal{M}_2 \cup r}^{\pi, h}(s) &= \sum_{a, s'} \pi_h(a|s) \left((P_1(s'|s, a) - P_2(s'|s, a)) V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s') + P_1(s'|s, a) (V_{\mathcal{M}_1 \cup r}^{\pi, h+1}(s') - V_{\mathcal{M}_2 \cup r}^{\pi, h+1}(s')) \right) \end{aligned}$$

Together, the first two results imply the third one because all terms in the sums are non-negative. \blacksquare

Lemma C.6 *Let $\mathcal{M}_1, \mathcal{M}_2$ be two MDP\(R) with transition dynamics P_1, P_2 respectively, r a reward function, and π_1^*, π_2^* optimal policy in $\mathcal{M}_1 \cup r$ and $\mathcal{M}_2 \cup r$, respectively. Then, for any state s and timestep h :*

$$\begin{aligned} V_{\mathcal{M}_1 \cup r}^{*, h}(s) - V_{\mathcal{M}_2 \cup r}^{*, h}(s) &\leq \sum_{h'=h} \sum_{s', a', s''} \eta_{\mathcal{M}_2, \pi_1^*}^{h, h'}(s'; s) \pi_{1, h}^*(a'|s') (P_1(s''|s', a') - P_2(s''|s', a')) V_{\mathcal{M}_1 \cup r}^{*, h}(s'') \\ V_{\mathcal{M}_2 \cup r}^{*, h}(s) - V_{\mathcal{M}_1 \cup r}^{*, h}(s) &\leq \sum_{h'=h} \sum_{s', a', s''} \eta_{\mathcal{M}_2, \pi_2^*}^{h, h'}(s'; s) \pi_{2, h}^*(a'|s') (P_2(s''|s', a') - P_1(s''|s', a')) V_{\mathcal{M}_2 \cup r}^{*, h}(s'') \end{aligned}$$

Proof

$$\begin{aligned} V_{\mathcal{M}_1 \cup r}^{*, h}(s) - V_{\mathcal{M}_2 \cup r}^{*, h}(s) &= \sum_{a, s'} \left(\pi_{1, h}^*(a|s) P_1(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') - \pi_{2, h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_2 \cup r}^{\pi_2^*, h+1}(s') \right. \\ &\quad \left. \pm \pi_{1, h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') \pm \pi_{1, h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_2 \cup r}^{\pi_2^*, h+1}(s') \right) \\ &= \sum_{a, s'} \left(\pi_{1, h}^*(a|s) P_2(s'|s, a) (V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') - V_{\mathcal{M}_2 \cup r}^{\pi_2^*, h+1}(s')) \right. \\ &\quad \left. + \pi_{1, h}^*(a|s) (P_1(s'|s, a) - P_2(s'|s, a)) V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') \right. \\ &\quad \left. + (\pi_{1, h}^*(a|s) - \pi_{2, h}^*(a|s)) P_2(s'|s, a) V_{\mathcal{M}_2 \cup r}^{\pi_2^*, h+1}(s') \right) \\ &\leq \sum_{a, s'} \left(\pi_{1, h}^*(a|s) P_2(s'|s, a) (V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') - V_{\mathcal{M}_2 \cup r}^{\pi_2^*, h+1}(s')) \right. \\ &\quad \left. + \pi_{1, h}^*(a|s) (P_1(s'|s, a) - P_2(s'|s, a)) V_{\mathcal{M}_1 \cup r}^{\pi_1^*, h+1}(s') \right) \end{aligned}$$

where the last inequality uses that π^* is optimal for $\mathcal{M}_2 \cup r$. Unrolling the recursion gives the first result. A similar argument yields the second results:

$$\begin{aligned}
V_{\mathcal{M}_2 \cup r}^{*,h}(s) - V_{\mathcal{M}_1 \cup r}^{*,h}(s) &= \sum_{a,s'} \left(\pi_{2,h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_2 \cup r}^{\pi_{2,h}^*,h+1}(s') - \pi_{1,h}^*(a|s) P_1(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s') \right. \\
&\quad \left. \pm \pi_{2,h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s') \right) \\
&= \sum_{a,s'} \left(\pi_{2,h}^*(a|s) P_2(s'|s, a) (V_{\mathcal{M}_2 \cup r}^{\pi_{2,h}^*,h+1}(s') - V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s')) \right. \\
&\quad \left. + \pi_{2,h}^*(a|s) P_2(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s') - \pi_{1,h}^*(a|s) P_1(s'|s, a) V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s') \right) \\
&\leq \sum_{a,s'} \left(\pi_{2,h}^*(a|s) P_2(s'|s, a) (V_{\mathcal{M}_2 \cup r}^{\pi_{2,h}^*,h+1}(s') - V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s')) \right. \\
&\quad \left. + \pi_{2,h}^*(a|s) (P_2(s'|s, a) - P_1(s'|s, a)) V_{\mathcal{M}_1 \cup r}^{\pi_{1,h}^*,h+1}(s') \right)
\end{aligned}$$

■

C.2 Feasible Reward Set

In this section, we characterize the feasible reward set first implicitly, then explicitly, and prove a result about error propagation. Metelli et al. [2021] provide a similar analysis in the infinite horizon setting.

Lemma C.7 (Feasible Reward Set Implicit) *A reward function r is feasible if and only if for all s, a, h it holds that: $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) = 0$ if $\pi_h^E(a|s) \geq 0$ and $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) \leq 0$ if $\pi_h^E(a|s) = 0$. Moreover, if the second inequality is strict, π^E is uniquely optimal, i.e., $\Pi_{\mathcal{M} \cup r}^* = \{\pi^E\}$.*

Proof The result follows directly from the definition of the feasible reward set. ■

Lemma C.8 *A Q -function satisfies the conditions of Theorem 2 if and only if there exists an $\{A_h \in \mathbb{R}_{\geq 0}^{S \times \mathcal{A}}\}_{h \in H}$ and $\{V_h \in \mathbb{R}^S\}$ such that for every $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$:*

$$Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = -A_h(s, a) \mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s)$$

Proof We first show that if $Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a)$ has this form, the conditions of Theorem 2 are satisfied, and then the converse. Assume $Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = -A_h(s, a) \mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s)$. Then,

$$V_{\mathcal{M} \cup r}^{\pi^E, h}(s) = \sum_a \pi_h^E(a|s) Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = V_h(s).$$

If $\pi_h^E(a|s) > 0$, then $Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = V_{\mathcal{M} \cup r}^{\pi^E, h}(s)$, which is the first condition of Theorem 2. If $\pi_h^E(a|s) = 0$, $Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = V_{\mathcal{M} \cup r}^{\pi^E, h}(s) - A_h(s, a) \leq V_{\mathcal{M} \cup r}^{\pi^E, h}(s)$, which is the second condition of Theorem 2.

For the converse, assume that the conditions of Theorem 2 hold, and let $V_h(s) = V_{\mathcal{M} \cup r}^{\pi^E, h}(s)$ and $A_h(s, a) = V_{\mathcal{M} \cup r}^{\pi^E, h}(s) - Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a)$. \blacksquare

Lemma C.9 (Feasible Reward Set Explicit) *A reward function r is feasible if and only if there exists an $\{A_h \in \mathbb{R}_{\geq 0}^{S \times A}\}_{h \in [H]}$ and $\{V_h \in \mathbb{R}^S\}_{h \in [H]}$ such that for all s, a, h it holds that:*

$$r_h(s, a) = -A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s) + \sum_{s'} P(s'|s, a)V_{h+1}(s')$$

Proof Since $Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) = r_h(s, a) + \sum_{s'} P(s'|s, a)V_{h+1}(s')$, using Theorem C.8, we have:

$$\begin{aligned} r_h(s, a) &= Q_{\mathcal{M} \cup r}^{\pi^E, h}(s, a) - \sum_{s'} P(s'|s, a)V_{h+1}(s') \\ &= -A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s) + \sum_{s'} P(s'|s, a)V_{h+1}(s') \end{aligned}$$

\blacksquare

Theorem 4 (Error Propagation) *Let (\mathcal{M}, π^E) and $(\widehat{\mathcal{M}}, \widehat{\pi}^E)$ be two IRL problems. Then, for any $r \in \mathcal{R}_{(\mathcal{M}, \pi^E)}$ there exists $\widehat{r} \in \widehat{\mathcal{R}}_{(\widehat{\mathcal{M}}, \widehat{\pi}^E)}$ such that:*

$$|r_h(s, a) - \widehat{r}_h(s, a)| \leq A_h(s, a)|\pi_h^E(a|s) - \widehat{\pi}_h^E(a|s)| + \sum_{s'} V_{h+1}(s')|P(s'|s, a) - \widehat{P}(s'|s, a)|$$

and we can bound $V_h \leq (H - h)R_{\max}$ and $A_h \leq (H - h)R_{\max}$.

Proof We start by rewriting r and \widehat{r} using Theorem 3:

$$\begin{aligned} r_h(s, a) &= -A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s) + \sum_{s'} P(s'|s, a)V_{h+1}(s') \\ \widehat{r}_h(s, a) &= -\widehat{A}_h(s, a)\mathbb{1}_{\{\widehat{\pi}_h^E(a|s)=0\}} + \widehat{V}_h(s) + \sum_{s'} \widehat{P}(s'|s, a)\widehat{V}_{h+1}(s') \end{aligned}$$

We can choose (w.l.o.g.) $V_h = \widehat{V}_h$ and $\widehat{A}_h = \mathbb{1}_{\{\pi_h^E(a|s)=0\}}A_h$:

$$\begin{aligned} r_h(s, a) - \widehat{r}_h(s, a) &= -A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}} + V_h(s) + \sum_{s'} P(s'|s, a)V_{h+1}(s') \\ &\quad + A_h(s, a)\mathbb{1}_{\{\widehat{\pi}_h^E(a|s)=0\}}\mathbb{1}_{\{\pi_h^E(a|s)=0\}} - V_h(s) - \sum_{s'} \widehat{P}(s'|s, a)V_{h+1}(s') \\ &= A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}}(\mathbb{1}_{\{\widehat{\pi}_h^E(a|s)=0\}} - 1) + \sum_{s'} V_{h+1}(s')(P(s'|s, a) - \widehat{P}(s'|s, a)) \\ &= -A_h(s, a)\mathbb{1}_{\{\pi_h^E(a|s)=0\}}\mathbb{1}_{\{\widehat{\pi}_h^E(a|s) \geq 0\}} + \sum_{s'} V_{h+1}(s')(P(s'|s, a) - \widehat{P}(s'|s, a)) \end{aligned}$$

The result follows by taking the absolute value and applying the triangle inequality. \blacksquare

Algorithm 2 Uniform sampling IRL with a generative model.

- 1: **Input:** significance $\delta \in (0, 1)$, target accuracy ϵ , maximum number of samples per iter. n_{\max}
 - 2: Initialize $k \leftarrow 0$, $\epsilon_0 \leftarrow H$
 - 3: **while** $\epsilon_k > \epsilon/2$ **do**
 - 4: Uniformly sample $\lceil \frac{n_{\max}}{SAH} \rceil$ samples from all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$
 - 5: For all samples, observe sample from transition dynamics and expert policy
 - 6: $k \leftarrow k + 1$
 - 7: Update \widehat{P}_k , $\widehat{\pi}_k$, and C_k^h
 - 8: Update accuracy $\epsilon_k \leftarrow H \max_{s,a,h} C_k^h(s, a)$
 - 9: **end while**
-

C.3 Uniform Sampling IRL with a Generative Model

In this section, we derive sample complexity results for uniform sampling with a generative model (Algorithm 2). Metelli et al. [2021] proved an analogous result for the infinite horizon setting focusing on transferable rewards. In contrast, our focus is on the finite horizon setting. Moreover, Metelli et al. [2021] considers to learn a reward that is transferable to a known target environment. In our setting, instead, we suppose to use the recovered reward function in the unknown source environment.

Definition C.10 (Optimality Criterion) *Let \mathcal{S} be a sampling strategy. Let $\mathcal{R}_{\mathfrak{B}}$ be the exact feasible set and $\mathcal{R}_{\widehat{\mathfrak{B}}}$ be the feasible set recovered after observing $n \geq 0$ samples collected from \mathcal{M} and π^E . We say that \mathcal{S} is (ϵ, δ, n) -correct if with probability at least $1 - \delta$ it holds that:*

$$\inf_{\widehat{r} \in \mathcal{R}_{\widehat{\mathfrak{B}}}} \sup_{\widehat{\pi}^* \in \Pi_{\widehat{\mathcal{M}} \cup \widehat{r}}^*} \max_{s,a,h} \left| Q_{\mathcal{M} \cup r}^{\pi^*,h}(s, a) - Q_{\mathcal{M} \cup r}^{\widehat{\pi}^*,h}(s, a) \right| \leq \epsilon \quad \text{for each } r \in \mathcal{R}_{\mathfrak{B}},$$
$$\inf_{r \in \mathcal{R}_{\mathfrak{B}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup r}^*} \max_{s,a,h} \left| Q_{\mathcal{M} \cup r}^{\pi^*,h}(s, a) - Q_{\mathcal{M} \cup r}^{\widehat{\pi}^*,h}(s, a) \right| \leq \epsilon \quad \text{for each } \widehat{r} \in \mathcal{R}_{\widehat{\mathfrak{B}}},$$

where π^* is an optimal policy in $\mathcal{M} \cup r$ and $\widehat{\pi}^*$ is an optimal policy in $\widehat{\mathcal{M}} \cup \widehat{r}$.

Lemma C.11 (Good Event) *Let π^E be a (possibly stochastic) expert policy. We estimate the expert policy with $\widehat{\pi}^E$ and the transition model P with an estimate \widehat{P}_k from k episodic interactions. Let $n_k^h(s, a)$ and $n_k^h(s)$ be the number of times state action pairs and states have been observed at time h within the first k episodes, and $n_k^{h+}(s, a) = \max\{1, n_k^h(s, a)\}$.*

Then,

$$\begin{aligned}
\mathbb{1}_{\{\pi_h^E(a|s)=0\}} \mathbb{1}_{\{\hat{\pi}_h^E(a|s) \geq 0\}} A_h(s, a) &\leq (H-h) R_{\max} \sqrt{\frac{\ell_k^h(s, a)}{n_k^{h+}(s, a)}} \\
\mathbb{1}_{\{\hat{\pi}_h^E(a|s)=0\}} \mathbb{1}_{\{\pi_h^E(a|s) \geq 0\}} \hat{A}_h(s, a) &\leq (H-h) R_{\max} \sqrt{\frac{\ell_k^h(s, a)}{n_k^{h+}(s, a)}} \\
\sum_{s'} |(P(s'|s, a) - \hat{P}_k(s'|s, a)) V_r^{\pi, h}(s')| &\leq (H-h) R_{\max} \sqrt{\frac{2\ell_k^h(s, a)}{n_k^{h+}(s, a)}} \\
\sum_{s'} |(P(s'|s, a) - \hat{P}_k(s'|s, a)) \hat{V}_r^{\pi, h}(s')| &\leq (H-h) R_{\max} \sqrt{\frac{2\ell_k^h(s, a)}{n_k^{h+}(s, a)}}
\end{aligned}$$

where $\ell_k^h(s, a) = \log\left(24SAH(n_k^{h+}(s, a))^2/\delta\right)$, holds simultaneously for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ and $k \geq 1$ with probability at least $1 - \delta$. We call the event that these equations hold the good event \mathcal{E} and write $P(\mathcal{E}) \geq 1 - \delta$.

Proof We show that each statement individually does not hold with probability less than $\delta/4$, which implies the result via a union bound. Let us denote $\beta_1(s, a, h) := (H-h) R_{\max} \sqrt{\frac{2\ell_k^h(s, a)}{n_k^{h+}(s, a)}}$. First, consider the last two inequalities. The probability that either of them does not hold is:

$$\begin{aligned}
&\Pr\left(\exists k \geq 1, (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : \sum_{s'} |(P(s'|s, a) - \hat{P}_k(s'|s, a)) V_r^{\pi, h}(s')| > \beta_1(s, a, h)\right) \\
&\stackrel{(a)}{\leq} \Pr\left(\exists m \geq 0, (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : \sum_{s'} |(P(s'|s, a) - \hat{P}_k(s'|s, a)) V_r^{\pi, h}(s')| > \beta_1(s, a, h)\right) \\
&\stackrel{(b)}{\leq} \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H \Pr\left(\sum_{s'} |(P(s'|s, a) - \hat{P}_k(s'|s, a)) V_r^{\pi, h}(s')| > \beta_1(s, a, h)\right) \\
&\stackrel{(c)}{\leq} \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H 2 \exp\left(-\frac{2\beta_1(s, a, h)^2 m^2}{4m(H-h)^2 R_{\max}^2}\right) \leq \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H 2 \exp(-\ell_k(s, a)) \\
&= \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H \frac{2\delta}{24SAH(m^+)^2} = \frac{\delta}{12} \left(1 + \sum_{m \geq 0} \frac{1}{m^2}\right) = \frac{\delta}{12} \left(1 + \frac{\pi^2}{6}\right) \leq \frac{\delta}{4}
\end{aligned}$$

Step (a) assumes that we visit a state action pair m times, and focuses on these m times the transition model for the given state-action pair is updated. Step (b) uses a union bound over m and (s, a) . Step (c) applies Hoeffding's inequality using that we estimate P with an average of samples, and $V_r^{\pi, h} \leq (H-h) R_{\max}$.

We show the first two inequalities similarly, with $\beta_2(s, a, h) := (H-h) R_{\max} \sqrt{\frac{\ell_k^h(s, a)}{n_k^{h+}(s, a)}}$

$$\begin{aligned}
& \Pr \left(\exists k \geq 1, (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : |(\pi_k^E(a|s) - \hat{\pi}_k^E(a|s))V_r^{\pi, h}(s')| > \beta_2(s, a, h) \right) \\
& \stackrel{(a)}{\leq} \Pr \left(\exists m \geq 0, (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : |(\pi_k^E(a|s) - \hat{\pi}_k^E(a|s))V_r^{\pi, h}(s')| > \beta_2(s, a, h) \right) \\
& \stackrel{(b)}{\leq} \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H \Pr \left(|(\pi_k^E(a|s) - \hat{\pi}_k^E(a|s))V_r^{\pi, h}(s')| > \beta_2(s, a, h) \right) \\
& \stackrel{(c)}{\leq} \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H 2 \exp \left(-\frac{2\beta_2(s, a, h)^2 m^2}{m(H-h)^2 R_{\max}^2} \right) \leq \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H 2 \exp(-\ell_k(s, a)) \\
& = \sum_{m \geq 0} \sum_{s, a} \sum_{h=0}^H \frac{2\delta}{24SAH(m^+)^2} = \frac{\delta}{12} \left(1 + \sum_{m \geq 0} \frac{1}{m^2} \right) = \frac{\delta}{12} \left(1 + \frac{\pi^2}{6} \right) \leq \frac{\delta}{4}
\end{aligned}$$

A union bound over all equations results in $P(\mathcal{E}) \geq 1 - \delta$. ■

Definition C.12 *We define the reward uncertainty as*

$$C_k^h(s, a) = (H - h)R_{\max} \min \left(1, 2\sqrt{\frac{2\ell_k^h(s, a)}{n_k^h(s, a)}} \right)$$

Corollary C.13 *Under the good event \mathcal{E} , in each iteration k it holds for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ that:*

$$|r_h(s, a) - \hat{r}_h^k(s, a)| \leq C_k^h(s, a)$$

Proof

$$\begin{aligned}
|r_h(s, a) - \hat{r}_h^k(s, a)| & \stackrel{(a)}{\leq} A_h(s, a) \mathbf{1}_{\{\pi_h^E(a|s)=0\}} \mathbf{1}_{\{\hat{\pi}_h^E(a|s) \geq 0\}} + \sum_{s'} V_{h+1}(s') |P(s'|s, a) - \hat{P}(s'|s, a)| \\
& \stackrel{(b)}{\leq} (H - h)R_{\max} \left(2\sqrt{\frac{2\ell_k^h(s, a)}{n_k^{h^+}(s, a)}} \right) = C_k^h(s, a)
\end{aligned}$$

where (a) uses Theorem 4 and (b) uses Theorem C.11. ■

Corollary C.14 *Let \mathcal{S} be a sampling strategy. Let $\mathcal{R}_{\mathfrak{B}}$ be the exact feasible set and $\mathcal{R}_{\hat{\mathfrak{B}}_k}$ be the feasible set recovered after k iterations. If*

$$H \max_{s, a, h} C_k^h(s, a) \leq \frac{\epsilon}{2},$$

then the conditions of Theorem 1 are satisfied.

Proof For the first condition of Theorem 1, observe:

$$\begin{aligned}
& \inf_{\hat{r} \in \mathcal{R}_{\mathfrak{B}_k}} \sup_{\hat{\pi}^* \in \Pi_{\mathcal{M} \cup \hat{r}}^*} \max_{s,a,h} (Q_{\mathcal{M} \cup r}^{\pi^*,h}(s,a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*,h}(s,a)) \\
& \stackrel{(a)}{\leq} \inf_{\hat{r} \in \mathcal{R}_{\mathfrak{B}_k}} \sup_{\hat{\pi}^* \in \Pi_{\mathcal{M} \cup \hat{r}}^*} \max_{s,a,h} \sum_{h'=h}^H \sum_{s',a'} \left(\eta_{\mathcal{M},\pi^*}^{h,h'}(s',a'|s,a) - \eta_{\mathcal{M},\hat{\pi}^*}^{h,h'}(s',a'|s,a) \right) (r_{h'}(s',a') - \hat{r}_{h'}(s',a')) \\
& \stackrel{(b)}{\leq} \inf_{\hat{r} \in \mathcal{R}_{\mathfrak{B}_k}} \sup_{\hat{\pi}^* \in \Pi_{\mathcal{M} \cup \hat{r}}^*} \max_{s,a,h} \left| \sum_{h'=h}^H \sum_{s',a'} \left(\eta_{\mathcal{M},\pi^*}^{h,h'}(s',a'|s,a) - \eta_{\mathcal{M},\hat{\pi}^*}^{h,h'}(s',a'|s,a) \right) C_k^{h'}(s',a') \right| \\
& \leq 2H \max_{s,a,h} C_k^h(s,a)
\end{aligned}$$

where (a) uses Theorem C.4 and (b) uses Theorem C.13.

For the second condition of Theorem 1, it follows similarly that:

$$\inf_{r \in \mathcal{R}_{\mathfrak{B}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup r}^*} \max_{s,a,h} (Q_{\mathcal{M} \cup r}^{\pi^*,h}(s,a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*,h}(s,a)) \leq 2H \max_{s,a,h} C_k^h(s,a)$$

Hence, if $H \max_{s,a,h} C_k^h(s,a) \leq \epsilon/2$, both conditions of Theorem 1 are satisfied. \blacksquare

Theorem C.15 (Sample Complexity of Uniform Sampling IRL) *With probability at least $1 - \delta$, Algorithm 2 stops at iteration τ fulfilling Theorem 1 with a number of samples upper bounded by:*

$$n \leq \tilde{O} \left(\frac{H^5 R_{\max}^2 SA}{\epsilon^2} \right)$$

Proof First, note

$$H \max_{s,a,h} C_k^h(s,a) = H^2 R_{\max} \max_{s,a,h} \left(2 \sqrt{\frac{2\ell_k^h(s,a)}{n_k^{h+}(s,a)}} \right)$$

After τ iterations, we have collected $\tau \cdot n_{\max}$ samples and for each s,a,h , we have:

$$n_{\tau}^{h+}(s,a) \geq \frac{\tau n_{\max}}{SAH} \geq 1$$

To terminate at iteration τ , we need to have for all s,a,h :

$$2H^2 R_{\max} \sqrt{\frac{2\ell_{\tau}^h(s,a)}{n_{\tau}^h(s,a)}} \leq \frac{\epsilon}{2}$$

which implies

$$n_{\tau}^h(s,a) \geq \frac{32H^4 R_{\max}^2 \ell_{\tau}^h(s,a)}{\epsilon^2}$$

By using Lemma B.8 by Metelli et al. [2021], we can conclude that the number of samples necessary to ensure accuracy ϵ is:

$$n \leq \tilde{O} \left(\frac{H^5 R_{\max}^2 SA}{\epsilon^2} \right)$$

■

Corollary C.16 *If the true reward function does not depend on the timestep h , i.e., $r_h(s, a) = r(s, a)$, then we can modify Algorithm 2 to only need $n \leq \tilde{\mathcal{O}}\left(\frac{H^4 R_{\max}^2 SA}{\epsilon^2}\right)$ samples.*

Proof If we know that the reward function does not depend on h we can choose $C_k(s, a) = \min_h C_k^h(s, a)$ as a confidence interval of the reward. Consequently, we can sample all states for a fixed h .

We still need for all s, a :

$$2H^2 R_{\max} \sqrt{\frac{2\ell_\tau^h(s, a)}{n_\tau^h(s, a)}} \leq \frac{\epsilon}{2} \Rightarrow n_\tau^h(s, a) \geq \frac{32H^4 R_{\max}^2 \ell_\tau^h(s, a)}{\epsilon^2}$$

Again, we use Lemma B.8 by Metelli et al. [2021], but we can eliminate one sum over H , ending up with:

$$n \leq \tilde{\mathcal{O}}\left(\frac{H^4 R_{\max}^2 SA}{\epsilon^2}\right)$$

■

C.4 Sample Complexity of AceIRL in Unknown Environments (Problem Independent)

We are now ready to analyze the sample complexity of AceIRL (Algorithm 1). We first consider the simple version of the algorithm: AceIRL Greedy. Then, we consider the full version of the algorithm after introducing a few additional lemma about the policy confidence set. We start by defining the error upper bound and deriving two lemmas that will help us to show that it is indeed an upper bound on the error we want to reduce.

Definition C.17 *We define recursively:*

$$E_k^H(s, a) = 0; \quad E_k^h(s, a) = \min\left((H - h)R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a' \in \mathcal{A}} E_k^{h+1}(s', a')\right)$$

where \hat{P} is the estimated transition model of the environment.

The first lemma shows that the error upper bound can upper bound the error due to estimating the transition model.

Lemma C.18 *Under the good event \mathcal{E} , for all policies π and reward functions r and all s, a, h :*

$$|Q_{\mathcal{M} \cup r}^{\pi, h}(s, a) - Q_{\mathcal{M} \cup r}^{\pi, h}(s, a)| \leq E_k^h(s, a)$$

Proof

$$\begin{aligned}
|Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h}(s, a)| &= \left| \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h+1}(s', a') \right. \\
&\quad \left. - \sum_{s'} P(s'|s, a) \sum_{a'} \pi(a'|s') Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s', a') \pm \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s', a') \right| \\
&\leq \left| \sum_{s'} (\widehat{P}(s'|s, a) - P(s'|s, a)) \sum_{a'} \pi(a'|s') Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s', a') \right| \\
&\quad + \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') \left| Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h+1}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s, a) \right| \\
&\leq C_k^h(s, a) + \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') \left| Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h+1}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s, a) \right|
\end{aligned}$$

For $h = H$ the result holds trivially. Now assuming it holds for $h + 1$, we consider step h :

$$\begin{aligned}
|Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h}(s, a)| &\leq C_k^h(s, a) + \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') \left| Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h+1}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s, a) \right| \\
&\leq C_k^h(s, a) + \sum_{s'} \widehat{P}(s'|s, a) \max_{a'} \left| Q_{\widehat{\mathcal{M}}_{\cup r}}^{\pi, h+1}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s, a) \right| \\
&\leq C_k^h(s, a) + \sum_{s'} \widehat{P}(s'|s, a) \max_{a'} E_k^{h+1}(s', a') = E_k^h(s, a)
\end{aligned}$$

■

The next lemma shows that the error upper bound can also upper bound the error in estimating the reward function, which is due to estimating the transition model and the expert policy.

Lemma C.19 *Under the good event \mathcal{E} , for all reward function r , all policies π , and all $s, a \in \mathcal{S} \times \mathcal{A}$:*

$$|Q_{\widehat{\mathcal{M}}_{\cup \hat{r}}}^{\pi, h}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h}(s, a)| \leq E_k^h(s, a)$$

Proof For $h = H$ the result holds trivially. Now assuming it holds for $h + 1$, we consider step h :

$$\begin{aligned}
&|Q_{\widehat{\mathcal{M}}_{\cup \hat{r}}}^{\pi, h}(s, a) - Q_{\mathcal{M}_{\cup r}}^{\pi, h}(s, a)| \\
&\leq |\hat{r}(s, a) - r(s, a)| + \sum_{s'} \widehat{P}(s'|s, a) \sum_{a'} \pi(a'|s') |Q_{\widehat{\mathcal{M}}_{\cup \hat{r}}}^{\pi, h+1}(s', a') - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s', a')| \\
&\leq |\hat{r}(s, a) - r(s, a)| + \sum_{s'} \widehat{P}(s'|s, a) \max_{a'} |Q_{\widehat{\mathcal{M}}_{\cup \hat{r}}}^{\pi, h+1}(s', a') - Q_{\mathcal{M}_{\cup r}}^{\pi, h+1}(s', a')| \\
&\leq |\hat{r}(s, a) - r(s, a)| + \sum_{s'} \widehat{P}(s'|s, a) \max_{a'} E_k^{h+1}(s', a') = E_k^h(s, a)
\end{aligned}$$

■

We can now combine the previous two lemmas to show that E is indeed an upper bound on the error we want to reduce. This implies correctness of AceIRL Greedy, which the following lemma formalizes.

Lemma C.20 (Correctness of AceIRL Greedy) *If AceIRL Greedy stops in episode k , after sampling n samples, i.e., $E_k^0(s_0, \pi_{k+1}(s_0)) \leq \frac{\epsilon}{4}$, then it fulfills Theorem 1.*

Proof Let us define the error

$$e_k^h(s, a) := |Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)|$$

where π^* is the true optimal policy in $\mathcal{M} \cup r$, and $\hat{\pi}^*$ is the optimal policy in $\widehat{\mathcal{M}} \cup \hat{r}$, i.e., in the estimated MDP using the inferred reward function. Then,

$$\begin{aligned} e_k^h(s, a) &= |Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a) \pm Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) \pm Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)| \\ &\leq \underbrace{|Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a)|}_{\leq E_k^h(s, a)} + \underbrace{|Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)|}_{\leq E_k^h(s, a)} \\ &\leq 2E_k^h(s, a) + |Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)| \end{aligned}$$

where, we used Theorem C.18.

Let us consider the remaining term $|Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)|$ in two steps. First, we have:

$$\begin{aligned} Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a) &\leq \underbrace{Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup \hat{r}}^{\pi^*, h}(s, a)}_{\leq E_k^h(s, a)} + \underbrace{Q_{\widehat{\mathcal{M}} \cup \hat{r}}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a)}_{\leq 0} \\ &\quad + \underbrace{Q_{\widehat{\mathcal{M}} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)}_{\leq E_k^h(s, a)} \leq 2E_k^h(s, a), \end{aligned}$$

where we used Theorem C.19 and the fact that $\hat{\pi}^*$ is optimal in the MDP $\widehat{\mathcal{M}} \cup \hat{r}$. Second, we have:

$$\begin{aligned} Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) &\leq \underbrace{Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)}_{\leq E_k^h(s, a)} + \underbrace{Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a)}_{\leq 0} \\ &\quad + \underbrace{Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a)}_{\leq E_k^h(s, a)} \leq 2E_k^h(s, a), \end{aligned}$$

where we used Theorem C.18 and the fact that π^* is optimal in the MDP $\mathcal{M} \cup r$. Overall, we find that

$$|Q_{\widehat{\mathcal{M}} \cup r}^{\pi^*, h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, h}(s, a)| \leq 2E_k^h(s, a),$$

and consequently,

$$e_k^h(s, a) \leq 4E_k^h(s, a).$$

Note that, $E_k^h(s, a)$ only sums positive terms, hence:

$$\max_{s,a,h} E_k^h(s, a) \leq \max_a E_k^0(s_0, a) = E_k^0(s_0, \pi_{k+1}(s_0))$$

Hence, if $E_k^0(s_0, \pi_{k+1}(s_0)) \leq \frac{\epsilon}{4}$, we have for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$:

$$e_k^h(s, a) \leq \epsilon$$

which implies correctness according to Theorem 1. ■

Next, we will analyze the sample complexity of AceIRL Greedy. Let us first define pseudo-counts that will be crucial to deal with the uncertainty of the transition dynamics in our analysis. This is similar to the analysis of UCRL for reward-free exploration by Kaufmann et al. [2021].

Definition C.21 *We define the pseudo-counts of visiting a specific state action pair at timestep h within the first k iterations as*

$$\bar{n}_k^h(s, a) := \sum_{i=1}^k \eta_{\mathcal{M}, \pi_i}^{0,h}(s, a | s_0),$$

where π_i is the exploration policy in episode i .

The following lemma allows us to introduce the pseudo-counts when considering the contraction of the reward confidence intervals.

Lemma C.22 *With probability at least $1 - \frac{\delta}{2}$ for all $s, a, h, k \in \mathcal{S} \times \mathcal{A} \times [H] \times \mathbb{N}^+$, we have:*

$$\min \left(\frac{2\ell_k^h(s, a)}{n_k^h(s, a)}, 1 \right) \leq \frac{8\bar{\ell}_k^h(s, a)}{\max(\bar{n}_k^h(s, a), 1)}$$

where $\bar{\ell}_k^h(s, a) = \log(24SAH(\bar{n}_k^h(s, a))^2/\delta)$.

Proof This result adapts Lemma 7 by Kaufmann et al. [2021] to our setting.

By Lemma 10 in Kaufmann et al. [2021], we have with probability at least $1 - \frac{\delta}{2}$:

$$n_k^h(s, a) \geq \frac{1}{2}\bar{n}_k^h(s, a) - \beta_{\text{cnt}}(\delta),$$

where $\beta_{\text{cnt}}(\delta) = \log(2SAH/\delta)$.

We distinguish two cases. First let $\beta_{\text{cnt}}(\delta) \leq \frac{1}{4}\bar{n}_k^h(s, a)$. Then $n_k^h(s, a) \geq \frac{1}{4}\bar{n}_k^h(s, a)$, and

$$\begin{aligned} \min \left(\frac{2\ell_k^h(s, a)}{n_k^h(s, a)}, 1 \right) &\leq \frac{2\ell_k^h(s, a)}{\max(n_k^h(s, a), 1)} = \frac{2 \log(24SAH(n_k^h(s, a))^2/\delta)}{\max(n_k^h(s, a), 1)} \\ &\leq \frac{2 \log(24SAH(\bar{n}_k^h(s, a)/4)^2/\delta)}{(\bar{n}_k^h(s, a)/4)} \leq \frac{8\bar{\ell}_k^h(s, a)}{\max(\bar{n}_k^h(s, a), 1)} \end{aligned}$$

where we use that $\log(24SAHx^2/\delta)/x$ is non-increasing for $x > 1$, and $\log(24SAHx^2/\delta)$ is non-decreasing and $\beta_{\text{cnt}}(\delta) \geq 1$.

Now consider let $\beta_{\text{cnt}}(\delta) > \frac{1}{4}\bar{n}_k^h(s, a)$. Then,

$$\min\left(\frac{2\ell_k^h(s, a)}{\bar{n}_k^h(s, a)}, 1\right) \leq 1 < 4\frac{\beta_{\text{cnt}}(\delta)}{\max(\bar{n}_k^h(s, a), 1)} \leq \frac{4\bar{\ell}_k^h(s, a)}{\max(\bar{n}_k^h(s, a), 1)}$$

where we used that $\ell_k^h(s, a) = \log(24SAH(\bar{n}_k^h(s, a))^2/\delta) = \beta_{\text{cnt}}(\delta) + \log(6\bar{n}_k^h(s, a))^2 \geq \beta_{\text{cnt}}(\delta)$. \blacksquare

The final lemma we need shows relates the error upper bound which is defined using our estimated transition model to a similar quantity defined using the (unknown) real transitions.

Lemma C.23 *Under the good event \mathcal{E} , we have for any s, a, h :*

$$E_k^h(s, a) \leq 2C_k^h(s, a) + \sum_{s'} P(s'|s, a) \max_{a'} E_k^{h+1}(s', a')$$

where P is the true transition model that we do not know.

Proof First note that $E_k^h(s, a) \leq H$ by definition. Now, consider:

$$\begin{aligned} E_k^h(s, a) &\leq C_k^h(s, a) + \sum_{s'} \widehat{P}(s'|s, a) \max_{a'} E_k^{h+1}(s', a') \\ &= C_k^h(s, a) + \sum_{s'} (\widehat{P}(s'|s, a) - P(s'|s, a) + P(s'|s, a)) \max_{a'} E_k^{h+1}(s', a') \\ &= C_k^h(s, a) + \underbrace{\sum_{s'} (\widehat{P}(s'|s, a) - P(s'|s, a)) \max_{a'} E_k^{h+1}(s', a')}_{\leq C_k^h(s, a)} + \sum_{s'} P(s'|s, a) \max_{a'} E_k^{h+1}(s', a') \\ &\leq 2C_k^h(s, a) + \sum_{s'} P(s'|s, a) \max_{a'} E_k^{h+1}(s', a') \end{aligned}$$

where we used the good event and the fact that C_k^h can only shrink over episodes. \blacksquare

Finally, we can analyze the sample complexity of AceIRL Greedy.

Theorem C.24 (AceIRL Greedy Sample Complexity (problem independent)) *AceIRL Greedy terminates with an (ϵ, δ, n) -correct solution, with*

$$n \leq \tilde{O}\left(\frac{H^5 R_{\max}^2 SA}{\epsilon^2}\right).$$

Proof Theorem C.20 shows that if AceIRL Greedy terminates, then it returns a (ϵ, δ, n) -correct solution. So, we need to show that it terminates within τ iterations and bound τ .

Let us consider the average error, defined by

$$\begin{aligned}
q_k^h &:= \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) E_k^h(s, a) \\
&\stackrel{(a)}{\leq} \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) (2C_k^h(s, a) + \sum_{s'} P(s'|s, a) \max_{a'} E_k^{h+1}(s', a')) \\
&= \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) (2C_k^h(s, a) + \sum_{s'} P(s'|s, a) \sum_{a'} \pi_{k+1}(a'|s') E_k^{h+1}(s', a')) \\
&= 2 \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) C_k^h(s, a) + q_k^{h+1}
\end{aligned}$$

where we used Theorem C.23 in step (a). Unrolling the recursion, results in:

$$q_k^h \leq 2 \sum_{h'=h}^H \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h'}(s, a|s_0) C_k^{h'}(s, a)$$

If the algorithm terminates at τ , we have for each $k < \tau$, and $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$: $\epsilon < 4E_k^0(s_0, \pi_{k+1}(s_0))$. We have $q_k^0 = E_k^0(s_0, \pi_{k+1}(s_0))$; therefore, as long we haven't stopped, we have $\epsilon \leq 4q_k^0$. Writing out this inequality, yields:

$$\begin{aligned}
\epsilon &\leq 4q_k^0 \leq 8 \sum_{h=0}^H \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) C_k^h(s, a) \\
&\leq 4HR_{\max} \sum_{h=0}^H \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) \sqrt{\frac{8 \log(12SAH(n_k^h(s, a))^2/\delta)}{\max(n_k^h(s, a), 1)}}
\end{aligned}$$

Using Theorem C.22, we can relate this to the pseudo-counts

$$\begin{aligned}
\epsilon &< 4HR_{\max} \sum_{h=0}^H \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) \sqrt{\frac{8 \log(12SAH(\bar{n}_k^h(s, a))^2/\delta)}{\max(\bar{n}_k^h(s, a), 1)}} \\
&\leq 4HR_{\max} \sum_{h=0}^H \sum_{s,a} \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) \sqrt{\frac{8 \log(12SAHk^2/\delta)}{\max(\bar{n}_k^h(s, a), 1)}}
\end{aligned}$$

Summing the inequality over $k = 0, \dots, T$ with $T < \tau$, we obtain

$$\begin{aligned}
\epsilon(T+1) &\leq 4HR_{\max} \sqrt{8 \log(12SAHT^2/\delta)} \sum_{h=0}^H \sum_{s,a} \sum_{k=1}^T \eta_{\mathcal{M}, \pi_{k+1}}^{0,h}(s, a|s_0) \frac{1}{\sqrt{\max(\bar{n}_k^h(s, a), 1)}} \\
&= 4HR_{\max} \sqrt{8 \log(12SAHT^2/\delta)} \sum_{h=0}^H \sum_{s,a} \sum_{k=1}^T \frac{\bar{n}_h^{k+1}(s, a) - \bar{n}_h^k(s, a)}{\sqrt{\max(\bar{n}_k^h(s, a), 1)}}
\end{aligned}$$

where we used the definition of the pseudo-counts in the last equality. Using Lemma 19 by Jaksch et al. [2010], we can further bound the sum in k :

$$\begin{aligned}\epsilon(T+1) &= 4HR_{\max}\sqrt{8\log(12SAHT^2/\delta)}\sum_{h=0}^H\sum_{s,a}\sqrt{\bar{n}_h^{T+1}(s,a)} \\ &\leq 4HR_{\max}\sqrt{8\log(12SAHT^2/\delta)}\sqrt{SA}\sum_{h=0}^H\sqrt{\sum_{s,a}\bar{n}_h^{T+1}(s,a)} \\ &= 4H^2R_{\max}\sqrt{8\log(12SAHT^2/\delta)}\sqrt{SA}\sqrt{T+1}\end{aligned}$$

It follows that

$$\begin{aligned}\epsilon\sqrt{T+1} &\leq 4H^2R_{\max}\sqrt{8SA\log(12SAHT^2/\delta)} \\ \epsilon^2\tau &\leq 128H^4R_{\max}^2SA\log(12SAH(\tau-1)^2/\delta)\end{aligned}$$

setting $\tau = T + 1$.

For large enough τ , this inequality cannot hold because $\sqrt{T+1}$ on the l.h.s grows faster than $\log(\tau)$ on the r.h.s. Hence, the stopping time τ is finite. Further, we can apply Lemma 15 by Kaufmann et al. [2021], and follow that

$$\tau \leq \tilde{O}\left(\frac{H^4R_{\max}^2SA}{\epsilon^2}\right)$$

If we observe H samples in each iteration, i.e., $N_E = 1$, we get a sample complexity of

$$n \leq \tilde{O}\left(\frac{H^5R_{\max}^2SA}{\epsilon^2}\right)$$

■

C.5 Sample Complexity of AceIRL in Unknown Environments (Problem Dependent)

For the problem dependent analysis, we will need this additional lemma also used by Kakade and Langford [2002].

Lemma C.25 (Lemma 6.1 by Kakade and Langford [2002]) *For any policy π :*

$$V_{\mathcal{M}\cup r}^{\pi^*,h}(s) - V_{\mathcal{M}\cup r}^{\pi,h}(s) = -\sum_{s',a'}\sum_{h'=h}^H\eta_{\mathcal{M},\pi}^{h,h'}(s',a';s)A_{\mathcal{M}\cup r}^{*,h'}(s',a')$$

Proof

$$\begin{aligned}
& V_{\mathcal{M} \cup r}^{*,h}(s) - V_{\mathcal{M} \cup r}^{\pi,h}(s) \\
&= \sum_a \pi_h^*(a|s) \left(r_h(s,a) + \sum_{s'} P(s'|s,a) V_{\mathcal{M} \cup r}^{*,h+1}(s') \right) \\
&\quad - \sum_a \pi_h(a|s) \left(r_h(s,a) + \sum_{s'} P(s'|s,a) V_{\mathcal{M} \cup r}^{\pi,h+1}(s') \right) \pm \sum_{a,s'} \pi_h(a|s) P(s'|s,a) V_{\mathcal{M} \cup r}^{*,h+1}(s') \\
&= \sum_a (\pi_h^*(a|s) - \pi_h(a|s)) r(s,a) + \sum_{a,s'} (\pi_h^*(a|s) - \pi_h(a|s)) P(s'|s,a) V_{\mathcal{M} \cup r}^{*,h+1}(s') \\
&\quad + \sum_{a,s'} \pi_h(a|s) P(s'|s,a) (V_{\mathcal{M} \cup r}^{*,h+1}(s) - V_{\mathcal{M} \cup r}^{\pi,h+1}(s)) \\
&= - \sum_a \pi(a|s) A_{\mathcal{M} \cup r}^{*,h}(s,a) + \sum_{a,s'} \pi_h(a|s) P(s'|s,a) (V_{\mathcal{M} \cup r}^{*,h+1}(s) - V_{\mathcal{M} \cup r}^{\pi,h+1}(s))
\end{aligned}$$

Unrolling the recursion yields the result. ■

We can now start with the analysis. First, we define the policy confidence set, and show that it indeed contains the relevant policies under the good event.

Definition C.26 We define the policy confidence set as

$$\hat{\Pi}_k = \{\pi | V_{\widehat{\mathcal{M}} \cup \hat{r}}^*(s_0) - V_{\widehat{\mathcal{M}} \cup \hat{r}}^\pi(s_0) \leq 10\epsilon_k\}$$

where $\hat{r} = \mathcal{A}(\mathcal{R}_{\hat{\mathfrak{B}}})$ is the reward estimated using an IRL algorithm \mathcal{A} . We choose ϵ_k recursively by solving the optimization problem

$$\epsilon_k = \max_{\pi \in \hat{\Pi}_{k-1}} \sum_{h=0}^H \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi}^{0,h}(s',a';s_0) C_k^h(s',a')$$

starting with $\epsilon_0 = \frac{1}{10}H$.

The following lemma will help us to deal with uncertainty about the transition dynamics.

Lemma C.27 Under the good event \mathcal{E} , if $\pi \in \hat{\Pi}_k$, then:

$$\begin{aligned}
|V_{\widehat{\mathcal{M}} \cup \hat{r}}^{\pi,h}(s) - V_{\mathcal{M} \cup r}^{\pi,h}(s)| &\leq \epsilon_k \\
|V_{\widehat{\mathcal{M}} \cup \hat{r}}^{*,h}(s) - V_{\mathcal{M} \cup r}^{*,h}(s)| &\leq \epsilon_k
\end{aligned}$$

Proof First by Theorem C.5:

$$\begin{aligned}
|V_{\widehat{\mathcal{M}} \cup \hat{r}}^{\pi,h}(s) - V_{\mathcal{M} \cup r}^{\pi,h}(s)| &\leq \sum_{h'=h}^H \sum_{s',a',s''} \eta_{\widehat{\mathcal{M}},\pi}^{h,h'}(s';s) \pi_{h'}(a'|s') |\widehat{P}(s''|s',a') - P(s''|s',a')| V_{\mathcal{M} \cup r}^{\pi,h'+1}(s'') \\
&\leq \sum_{h'=h}^H \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi}^{h,h'}(s';s) \pi_{h'}(a'|s') C_k(s',a') \leq \epsilon_k
\end{aligned}$$

Then, by Theorem C.6:

$$\begin{aligned} V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) &\leq \sum_{h'=h} \sum_{s',a',s''} \eta_{\widehat{\mathcal{M}},\pi^*}^{h,h'}(s';s)\pi_{h'}^*(a'|s')(P(s''|s',a') - \widehat{P}(s''|s',a'))V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s'') \\ &\leq \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi^*}^{h,h'}(s';s)\pi_{h'}^*(a'|s')C_k(s',a') \leq \epsilon_k \end{aligned}$$

And, similarly

$$\begin{aligned} V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) &\leq \sum_{h'=h} \sum_{s',a',s''} \eta_{\widehat{\mathcal{M}},\widehat{\pi}^*}^{h,h'}(s';s)\widehat{\pi}_{h'}^*(a'|s')(\widehat{P}(s''|s',a') - P(s''|s',a'))V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s'') \\ &\leq \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\widehat{\pi}^*}^{h,h'}(s';s)\widehat{\pi}_{h'}^*(a'|s')C_k(s',a') \leq \epsilon_k \end{aligned}$$

■

Now we show that the relevant policies are always in the policy confidence set, conditioned on the good event.

Lemma C.28 *Conditioned the good event \mathcal{E} , if $\pi^*, \widehat{\pi}^* \in \widehat{\Pi}_{k-1}$, then $\pi^* \in \widehat{\Pi}_k$.*

Proof Let $r \in \mathcal{R}_{\mathfrak{B}}$. Then

$$\begin{aligned} V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi^*,h}(s) &= V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) + V_{\widehat{\mathcal{M}}\cup r}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi^*,h}(s) \\ &\stackrel{(a)}{\leq} \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi^*}^{h,h'}(s',a'|s)C_k^{h'}(s',a') + \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi^*}^{h,h'}(s',a'|s)C_k^{h'}(s',a') \stackrel{(b)}{\leq} 2\epsilon_k \end{aligned}$$

where (a) uses Theorem C.2, Theorem C.3 and Theorem C.13, (b) uses that $\pi^* \in \widehat{\Pi}_{k-1}$ and the definition of ϵ_k . Hence,

$$\max_s \left(V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{*,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi^*,h}(s) \right) \leq 2\epsilon_k \leq 10\epsilon_k$$

and therefore $\pi^* \in \widehat{\Pi}_k$. ■

Lemma C.29 *Conditioned on the good event \mathcal{E} , for every policy π and episodes $k' > k$, there exists $\widehat{r}_{k'} \in \mathcal{R}_{\mathfrak{B}_{k'}}$, such that:*

$$\max_s \left(V_{\widehat{\mathcal{M}}\cup \widehat{r}_{k'}}^{\pi,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi,h}(s) \right) \leq 4\epsilon_k$$

Proof Similarly to the proof of the previous lemma, we have

$$\begin{aligned} V_{\widehat{\mathcal{M}}\cup \widehat{r}_{k'}}^{\pi,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi,h}(s) &= V_{\widehat{\mathcal{M}}\cup \widehat{r}_{k'}}^{\pi,h}(s) - V_{\widehat{\mathcal{M}}\cup r}^{\pi,h}(s) + V_{\widehat{\mathcal{M}}\cup r}^{\pi,h}(s) - V_{\widehat{\mathcal{M}}\cup \widehat{r}_k}^{\pi,h}(s) \\ &\leq \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi}^{h,h'}(s',a'|s)C_{k'}^{h'}(s',a') + \sum_{h'=h} \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\pi}^{h,h'}(s',a'|s)C_k^{h'}(s',a') \leq 2\epsilon_k \end{aligned}$$

where we use that the confidence intervals are shrinking with increasing episode number, i.e., $\epsilon_{k'} \leq \epsilon_k$.

By combining this with Theorem C.27, we get the result:

$$\begin{aligned} & \max_s \left(V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s) \right) \\ &= \max_s \left(\underbrace{V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s) - V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s)}_{\leq \epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s) - V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s)}_{\leq 2\epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s)}_{\leq \epsilon_k} \right) \leq 4\epsilon_k \end{aligned}$$

■

Lemma C.30 *Under the good event \mathcal{E} , if $\hat{\pi}_{k'}^*$, $\pi \in \hat{\Pi}_{k-1}$ and $\pi \notin \hat{\Pi}_k$, then the policy π is suboptimal for some reward $\hat{r}_{k'} \in \mathcal{R}_{\hat{\mathfrak{B}}_{k'}}$ for all $k' \geq k$.*

Proof We can observe that

$$\begin{aligned} & V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s_0) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{*, h}(s_0) = V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s_0) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\hat{\pi}_{k'}^*, h}(s_0) \\ &= \underbrace{V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\pi, h}(s_0) - V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s_0)}_{\stackrel{(a)}{\leq} 4\epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s_0) - V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s_0)}_{\stackrel{(b)}{\leq} \epsilon_k} \\ &+ \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\pi, h}(s_0) - V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\hat{\pi}_{k'}^*, h}(s_0)}_{\stackrel{(c)}{>} 10\epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{\mathcal{R}}_k}^{\hat{\pi}_{k'}^*, h}(s_0) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_k}^{\hat{\pi}_{k'}^*, h}(s_0)}_{\stackrel{(b)}{\leq} \epsilon_k} \\ &+ \underbrace{V_{\mathcal{M} \cup \hat{\mathcal{R}}_k}^{\hat{\pi}_{k'}^*, h}(s_0) - V_{\mathcal{M} \cup \hat{\mathcal{R}}_{k'}}^{\hat{\pi}_{k'}^*, h}(s_0)}_{\stackrel{(a)}{\leq} 4\epsilon_k} > 0 \end{aligned}$$

where we applied (a) Theorem C.27, (b) Theorem C.29, and (c) the definition of $\hat{\Pi}_k$ and the fact that $\pi \notin \hat{\Pi}_k$. Consequently, π is suboptimal for at least some reward function $\hat{r}_{k'} \in \mathcal{R}_{\hat{\mathfrak{B}}_{k'}}$. ■

Corollary C.31 *For $\epsilon_0 = \frac{H}{10}$, for every $k \geq 0$ it holds that both π^* , $\hat{\pi}_{k+1}^* \in \hat{\Pi}_k$.*

Proof We show the statement by induction over k . For $k = 0$, we have $10\epsilon_0 = H$ and therefore $\hat{\Pi}_0$ contains all policies. Assume that for $k - 1$ the statement holds, i.e., π^* , $\hat{\pi}_k^* \in \hat{\Pi}_{k-1}$, and consider k . By Theorem C.28, $\pi^* \in \hat{\Pi}_k$. Note, that $\hat{\pi}_{k+1}^* \in \hat{\Pi}_{k-1}$. Hence, by Theorem C.29, it follows that $\hat{\pi}_{k+1}^* \in \hat{\Pi}_k$ because it would be suboptimal otherwise which is a contradiction. ■

The last result we need, is quantifying the size of the policy confidence set.

Lemma C.32 Under the good event \mathcal{E} , let $\tilde{r} \in \operatorname{argmin}_{r \in \mathcal{R}_{\mathfrak{B}}} \max_{s,a} (r(s,a) - \hat{r}_k(s,a))$, where $\hat{r}_k = \mathcal{A}(\mathcal{R}_{\mathfrak{B}_k})$. If $\pi \in \hat{\Pi}_k$, then $\max_s (V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{*,h}(s) - V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{\pi,h}(s)) \leq 12\epsilon_k$.

Proof

$$V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{*,h}(s) - V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{\pi,h}(s) = \underbrace{V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{*,h}(s) - V_{\widehat{\mathcal{M}} \cup \hat{r}_k}^{*,h}(s)}_{\leq \epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{r}_k}^{*,h}(s) - V_{\widehat{\mathcal{M}} \cup \hat{r}_k}^{\pi,h}(s)}_{\leq 10\epsilon_k} + \underbrace{V_{\widehat{\mathcal{M}} \cup \hat{r}_k}^{\pi,h}(s) - V_{\widehat{\mathcal{M}} \cup \tilde{r}}^{\pi,h}(s)}_{\leq \epsilon_k} \epsilon_k \leq 14\epsilon_k$$

■

Next, we define the error upper bound based on the policy confidence set.

Definition C.33 Using $\hat{\Pi}_k$, we define recursively:

$$\begin{aligned} \hat{E}_k^H(s, a) &= 0 \\ \hat{E}_k^h(s, a) &= \min\left((H - h)R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{\pi \in \hat{\Pi}_{k-1}} \pi(a'|s') \hat{E}_k^{h+1}(s', a')\right) \end{aligned}$$

where \hat{P} is the estimated transition model of the environment. In contrast to Theorem C.17, the maximization is over policies in $\hat{\Pi}_k$ rather than all actions.

This definition allows us to derive results that are analogous to the problem independent case.

Lemma C.34 Under the good event \mathcal{E} , for all policies $\pi \in \hat{\Pi}_k$ and reward functions r and all $s, a \in \mathcal{S} \times \mathcal{A}$:

$$|Q_{\widehat{\mathcal{M}} \cup r}^{\pi,h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\pi,h}(s, a)| \leq \hat{E}_k^h(s, a)$$

Proof The proof is the same as for Theorem C.18, restricting the set of policies to $\hat{\Pi}_k$. ■

Lemma C.35 Under the good event \mathcal{E} , for all reward function r , all policies $\pi \in \hat{\Pi}_k$, and all $s, a \in \mathcal{S} \times \mathcal{A}$:

$$|Q_{\widehat{\mathcal{M}} \cup \tilde{r}}^{\pi,h}(s, a) - Q_{\widehat{\mathcal{M}} \cup r}^{\pi,h}(s, a)| \leq \hat{E}_k^h(s, a)$$

Proof The proof is the same as for Theorem C.19, restricting the set of policies to $\hat{\Pi}_k$. ■

Lemma C.36 Under the good event \mathcal{E} , we have for any s, a, h :

$$\hat{E}_k^h(s, a) \leq 2C_k^h(s, a) + \sum_{s'} P(s'|s, a) \max_{\pi \in \hat{\Pi}_{k-1}} \pi(a'|s') \hat{E}_k^{h+1}(s', a')$$

Proof The proof is the same as for Theorem C.36. ■

Finally, we can combine these results to analyze the algorithm's sample complexity.

Theorem 5 [AceIRL Sample Complexity] AceIRL returns a (ϵ, δ, n) -correct solution with

$$n \leq \tilde{\mathcal{O}} \left(\min \left[\frac{H^5 R_{\max}^2 S A}{\epsilon^2}, \frac{H^4 R_{\max}^2 S A \epsilon_{\tau-1}^2}{\min_{s,a,h} (A_{\mathcal{M} \cup r}^{*,h}(s,a))^2 \epsilon^2} \right] \right)$$

where $\epsilon_{\tau-1}$ depends on the choice of N_E , the number of episodes of exploration in each iteration. $A_{\mathcal{M} \cup r}^{*,h}(s,a)$ is the advantage function of $r \in \operatorname{argmin}_{r \in \mathcal{R}_{\mathfrak{B}}} \max_{h,s,a} (r_h(s,a) - \hat{r}_{k,h}(s,a))$, the reward function from the feasible set $\mathcal{R}_{\mathfrak{B}}$ closest to the estimated reward function \hat{r}_k .

Proof First note that the analysis of Theorem C.24 still applies; so, in the worst case we get the same sample complexity. The key difference is that we no longer use the overall greedy policy w.r.t E_k^h , but restrict ourselves to policies in $\hat{\Pi}_k$.

Again, we consider the error

$$e_k^{\pi,h}(s,a) := |Q_{\mathcal{M} \cup r}^{\pi*,h}(s,a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}*,h}(s,a)|$$

where π^* is the true optimal policy in $\mathcal{M} \cup r$, and $\hat{\pi}^*$ is the optimal policy in $\widehat{\mathcal{M}} \cup \hat{r}$, i.e., in the estimated MDP using the inferred reward function.

Similar, to the proof of Theorem C.20, we can use Theorem C.34 and Theorem C.35 to show for all policies $\pi \in \hat{\Pi}_k^h$, that:

$$e_k^{\pi,h}(s,a) \leq 4\hat{E}_k^h(s,a)$$

which implies the correctness of the algorithm according to Theorem C.14 when stopping at

$$\hat{E}_k^0(s_0, \pi_{k+1}(s_0)) \leq \frac{\epsilon}{4} \quad (1)$$

Now, consider the following condition for all s, a, h :

$$C_k^h(s,a) \leq -A_{\mathcal{M} \cup \tilde{r}}^{*,h}(s,a) \frac{\epsilon}{48\epsilon_{k-1}}, \quad (2)$$

where $\tilde{r} \in \operatorname{argmin}_{r \in \mathcal{R}_{\mathfrak{B}}} \max_{h,s,a} (r_h(s,a) - \hat{r}_{k,h}(s,a))$. We will (a) show that when this condition holds the previous stopping condition also holds, and (b) analyze after how many iterations this condition will certainly hold. Together this will yield the result.

To show that Equation (2) implies Equation (1), we assume that Equation (2) holds. Then, we get by applying Theorem C.36 recursively:

$$\begin{aligned} \hat{E}_k^0(s_0, \pi_{k+1}(s_0)) &\leq 2 \max_{\pi \in \hat{\Pi}_{k-1}} \max_a \sum_{h=0}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{0,h}(s',a'; s_0, a) C_k^h(s',a') \\ &\leq 2 \max_{\pi \in \hat{\Pi}_{k-1}} \max_a \sum_{h=0}^H \sum_{s',a'} \eta_{\mathcal{M},\pi}^{0,h}(s',a'; s_0, a) \left(-A_{\mathcal{M} \cup \tilde{r}}^{*,h}(s',a') \frac{\epsilon}{48\epsilon_{k-1}} \right) \\ &\stackrel{(a)}{\leq} 2 \max_{\pi \in \hat{\Pi}_{k-1}} (V_{\mathcal{M} \cup r}^{*,0}(s_0) - V_{\mathcal{M} \cup r}^{\pi,0}(s_0)) \frac{\epsilon}{48\epsilon_{k-1}} \stackrel{(b)}{\leq} \frac{\epsilon}{4} \end{aligned}$$

where (a) uses Theorem C.25 and (b) uses Theorem C.32.

Next, we analyze after how many iterations Equation (2) holds, which will give a lower bound on the sample complexity result. The argument proceeds similar to the proof of Theorem C.24.

Before the algorithm terminates at τ , we have for all $k < \tau$:

$$\min_{s,a,h}(-A_{\widehat{\mathcal{M}}\cup\widehat{r}}^{*,h}(s,a))\frac{\epsilon}{48\epsilon_{k-1}} < \max_{s,a,h}C_k^h(s,a) \leq HR_{\max}\sqrt{\frac{2\ell_k^h(s,a)}{\max(N_k^h(s,a),)}}$$

Using similar argument to the proof of Theorem C.24, using the same pseudo-counts, we arrive at:

$$\min_{s,a,h}(-A_{\widehat{\mathcal{M}}\cup\widehat{r}}^{*,h}(s,a))\frac{\epsilon}{48\epsilon_{\tau-1}}\sqrt{\tau+1} \leq HR_{\max}\sqrt{8SA\log(12SAH\tau^2/\delta)}$$

Again, we can use Lemma 15 by Kaufmann et al. [2021] to find that

$$\tau \leq \tilde{\mathcal{O}}\left(\frac{H^3R_{\max}^2SA\epsilon_{\tau-1}^2}{\min_{s,a,h}(A_{\widehat{\mathcal{M}}\cup\widehat{r}}^{*,h}(s,a))^2\epsilon^2}\right)$$

■

C.6 Computing the Exploration Policy

To run AceIRL, we need to solve the optimization problem:

$$\pi_k^h = \min_{\pi} \max_{\hat{\pi} \in \widehat{\Pi}_{k-1}} \sum_{h=0}^H \sum_{s',a'} \eta_{\widehat{\mathcal{M}},\hat{\pi}}^{0,h}(s',a';s_0)\widehat{C}_k^h(s',a'|\pi)$$

For simplicity let us denote the state visitation frequencies by

$$\begin{aligned} \mu_h(s,a) &:= \eta_{\widehat{\mathcal{M}},\pi}^{0,h}(s,a;s_0) \\ \hat{\mu}_h(s,a) &:= \eta_{\widehat{\mathcal{M}},\hat{\pi}}^{0,h}(s,a;s_0) \end{aligned}$$

Let us introduce the following matrix notation

$$\tilde{A} = \begin{bmatrix} I & 0 & 0 & 0 & \dots & 0 \\ \widehat{P} & -I & 0 & 0 & \dots & 0 \\ 0 & \widehat{P} & -I & 0 & \dots & 0 \\ & & \dots & & & \\ 0 & 0 & \dots & 0 & \widehat{P} & -I \\ I & 0 & 0 & \dots & 0 & 0 \\ 0 & I & 0 & \dots & 0 & 0 \\ & & \dots & & & \\ 0 & 0 & 0 & \dots & I & 0 \\ 0 & 0 & 0 & \dots & 0 & I \end{bmatrix}, \quad a = \begin{bmatrix} \hat{r}_{k-1}^0 \\ \hat{r}_{k-1}^1 \\ \dots \\ \hat{r}_{k-1}^H \end{bmatrix}, \quad A = \begin{bmatrix} A & 0 \\ a^T & -1 \end{bmatrix},$$

$$x = \begin{bmatrix} \mu_0 \\ \mu_1 \\ \dots \\ \mu_H \\ t \end{bmatrix}, \quad \hat{x} = \begin{bmatrix} \hat{\mu}_0 \\ \hat{\mu}_1 \\ \dots \\ \hat{\mu}_H \end{bmatrix}, \quad b = \begin{bmatrix} \bar{\mu}_0 \\ 0 \\ \dots \\ 0 \\ 1 \\ \dots \\ 1 \\ -10\epsilon_{k-1} \end{bmatrix}, \quad c = \begin{bmatrix} C_0 \\ C_1 \\ \dots \\ C_H \\ 1 \end{bmatrix},$$

where $\bar{\mu}_0$ is the actual initial state distribution of the environment (which we assume to know). We can now write the inner maximization problem above as a linear program:

$$\begin{aligned} \max_x \quad & c^T x \\ & Ax = b \\ & x \geq 0 \end{aligned}$$

The corresponding dual problem is:

$$\begin{aligned} \min_y \quad & b^T y \\ & A^T y \geq c \end{aligned}$$

Using this we can write the full min-max problem as:

$$\begin{aligned} \min_{\hat{x}, y} \quad & b^T y \\ & A^T y \geq c(x) \\ & \tilde{A}x = b \\ & x \geq 0 \end{aligned}$$

which is a convex optimization problem, if we use:

$$C_h(s, a) = 2(H - h)R_{\max} \sqrt{\frac{2 \log(24SAH(\max(1, n_k^h(s, a)))^2 / \delta)}{\max(1, \hat{n}_{k+1}^h(s, a))}}$$

where $\hat{n}_{k+1}^h(s, a) = n_k^h(s, a) + \mu_h(s, a) * N_E$ is the number of times we expect h, s, a to be visited at the next iteration.

Solving this optimization problem yields the state-visitation frequencies $\hat{\mu}_k(s, a)$. We can then find the exploration policy that induces these state-visitations simply as:

$$\pi_{k,h}(a|s) := \frac{\hat{\mu}_k^h(s, a)}{\sum_{a'} \hat{\mu}_k^h(s, a')}.$$

Appendix D. Experimental Details

In this section, we provide more details on our experiments. We discuss the environments in detail (Appendix D.1), provide some information on the implementation and the libraries and computational resources we used (Appendix D.2), and we provide more full plots of all experiments we discussed (Appendix D.3).

D.1 Details on the Environments

Four Paths. The four paths environment has 41 states and 4 actions:

$$\mathcal{S} = \{c, l_1, \dots, l_{10}, u_1, \dots, u_{10}, r_1, \dots, r_{10}, d_1, \dots, d_{10}\}, \quad \mathcal{A} = \{a_1, a_2, a_3, a_4\},$$

and a time horizon of $H = 20$. The agent starts in the center state c , from which can move in four directions: left (a_1), up (a_2), right (a_3), or down (a_4). Each action a_i has a probability p_i of failing. If an action fails it moves in the opposite direction. p_1, \dots, p_4 are sampled uniformly from $(0, 0.3)$. One of the states $(l_{10}, u_{10}, r_{10}, d_{10})$ is chosen as the goal state at random. The reward in the goal state is 1, all other rewards are 0.

Double Chain. The *Double Chain* MDP, proposed by Kaufmann et al. [2021], consists of L states $\mathcal{S} = \{s_0, \dots, s_{L-1}\}$, and two actions $\mathcal{A} = \{\text{left}, \text{right}\}$, which correspond to a transition to the left or to the right. When the agent takes an action, there is a 0.1 probability of moving to the other direction. The state s_{L-1} has reward 1, all other states have reward 0, and the agent starts in the center of the chain at $s_{(L-1)/2}$. We choose $L = 31$, similar to Kaufmann et al. [2021]. The environment has horizon $H = 20$.

Chain. The *Chain* MDP, proposed by Metelli et al. [2021] has 6 states $\mathcal{S} = \{s_1, s_2, s_3, s_4, s_5, s_u\}$ and 10 actions $\mathcal{A} = \{a_1, \dots, a_{10}\}$. The agent starts in a random initial state. Taking action a_{10} moves it right along the chain with probability 0.7 and to state s_u with probability 0.3. Any other action moves the agent right with probability 0.3 and to state s_u with probability 0.7. If the agent is in state s_u , action a_{10} moves it back to state s_1 with probability 0.05. Any other action moves it to s_1 with probability 0.01. The reward is 1 in all states except s_u where the reward is 0. Metelli et al. [2021] provide an illustration of the environment in Figure 3. We choose $H = 10$ for the chain.

Gridworld. The *Gridworld*, proposed by Metelli et al. [2021], is a 3×3 gridworld with an obstacle in the center cell $(2, 2)$ and a goal cell at the right center cell $(2, 1)$. The agent starts in a random non-goal cell, and it has 4 action one to move in each direction. If the agent takes an action with probability 0.3 the action fails and the agent moves in a random direction instead. If the agent is in the center cell $(2, 2)$ which has the obstacle, if the agent would move right it instead stays in the center cell with probability 0.8. The reward in the goal cell is 1, all other rewards are 0. Metelli et al. [2021] provide an illustration of the gridworld in Figure 6. We choose $H = 10$ for the gridworld.

Random MDPs. We generate random MDPs by uniformly sampling an initial state distribution and transition matrix and normalizing them. The rewards are sampled uniformly between 0 and 1. Our random MDPs have 9 states, 4 actions and horizon 10.

D.2 Implementation Details

We provide a full implementation of AceIRL in Python, using multiple open sources libraries, including `cvxpy` and the SCS optimizer [Diamond and Boyd, 2016, O’Donoghue et al., 2016] for solving the optimization problem in Appendix C.6, and standard libraries for numerical computing, including `numpy`, and `scipy`. We choose Maximum Entropy IRL [Ziebart et al., 2008] as an IRL algorithm, but AceIRL is agnostic to this choice.

We ran experiments in parallel on a server with two 64 Core AMD EPYC 7742 2.25GHz processors. We estimate a total wall-clock time of less than 48 hours for running all experiments presented in this paper, including 50 random seeds each.

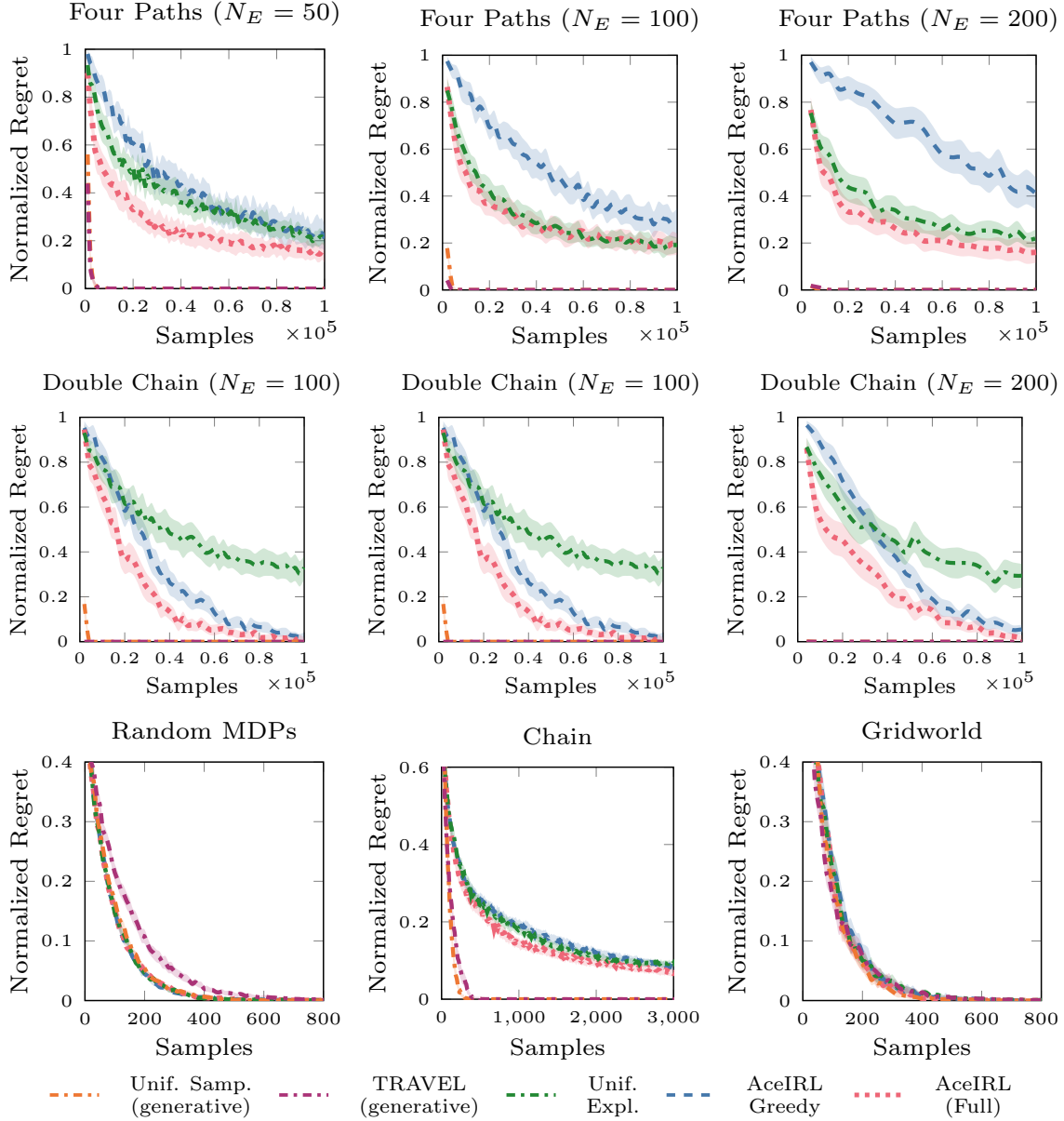


Figure D.2: Full learning curves for all experiments shown in Table A.1. Similar to Figure A.1, we show the mean and 95% confidence intervals computed over 50 random seeds. In addition to the exploration algorithms, we also show uniform sampling and TRAVEL which are much faster in most cases because they have access to a generative model.

D.3 Additional Results

We provide full learning curves for all experiments discussed in Figure D.2.

Appendix E. Connection to Reward-free Exploration

In the *reward-free exploration* problem, introduced by Jin et al. [2020], the agent explores an MDP \mathcal{R} to learn a transition model. In each iteration it chooses a new exploration policy based on previous data. The goal is to ensure that if the agent is given a reward function r after the exploration phase it can find a good policy using its transition model. Jin et al. [2020] formalize this goal as reducing the error:

$$V_{\mathcal{M} \cup r}^{\pi^*, 0}(s_0) - V_{\widehat{\mathcal{M}} \cup r}^{\hat{\pi}^*, 0},$$

where $\hat{\pi}^*$ is the optimal policy in the estimated MDP $\widehat{\mathcal{M}} \cup r$. Note the striking similarity between this problem, and the active IRL problem, we study in this paper. We want to reduce a similar error (cf. Theorem 1), but we have additional information about the reward in form of the expert policy.

The *Reward-free UCRL* algorithm, proposed by Kaufmann et al. [2021], is essentially analogous to AceIRL Greedy (Section 4.1). Reward-free UCRL explores greedily with respect to an upper bound on the value function error. However, the exploration policy needs to be updated after each episode to adapt to the new uncertainty estimates. This might be expensive or not possible in practice. Instead, we could consider a *batched* version of reward-free exploration, where in each iteration the agent explores for N_E episodes, similar to our Active IRL problem. In this setting, a greedy policy w.r.t. uncertainty is suboptimal because it does not adapt to the reduced uncertainty over the N_E episodes.

Instead, we can consider reducing the expected uncertainty at the next iteration, similar to our discussion in Section 4.2. If our error estimate is denoted by $E_k(s, a)$, we do no longer act greedily w.r.t. E_k . Instead we try to estimate the error at the next iteration $\hat{E}_{k+1}(s, a | \pi)$ as a function of the policy and try to select the policy that reduces this error. In the tabular case, we can formulate this as a convex optimization problem, analogous to Appendix C.6. We call this adaptation of AceIRL to the reward-free exploration problem *Ace-RF*.

Figure E.3 shows illustrative results of this algorithm in the batched reward-free exploration setting in the *Double Chain* environment. We find that for larger batch sizes, choosing an exploration policy that reduces future uncertainty is significantly better than reward-free UCRL.

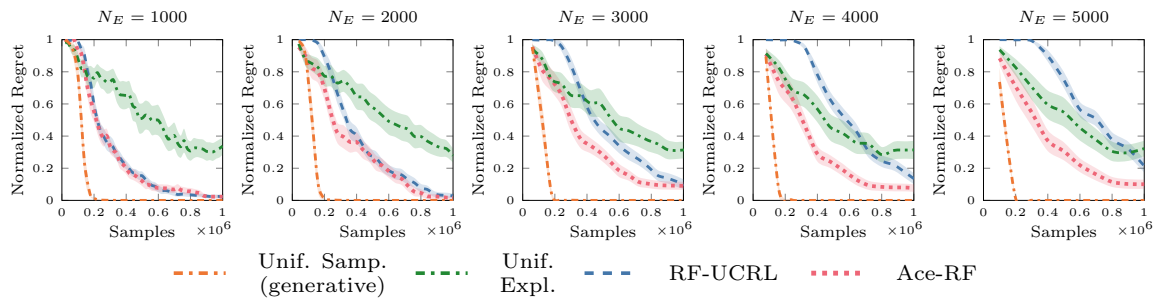


Figure E.3: Illustrative experiments for reward-free exploration in the *Double Chain* environment proposed by Kaufmann et al. [2021]. The difference to our Active IRL setting is that the agent does not have access to the expert policy during exploration, but still tries to learn a good model of the environment. During testing it then gets access to the reward function, and the regret measures the suboptimality of the policy trained in the agent’s transition model. We find that the ideas used in AceIRL are also useful for batched reward-free exploration with target N_E .