



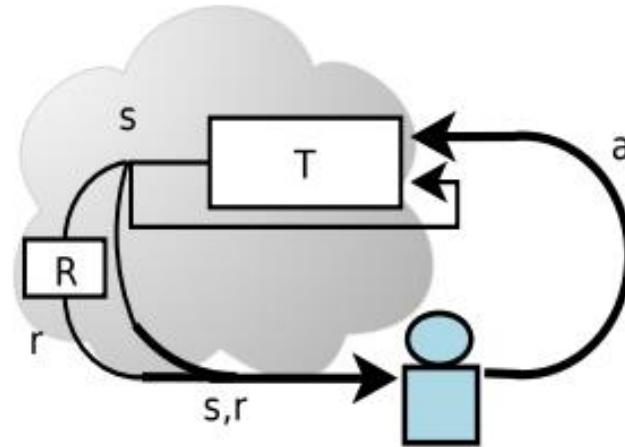
Active Exploration for Inverse Reinforcement Learning

David Lindner
Department of Computer Science
ETH Zurich
david.lindner@inf.ethz.ch

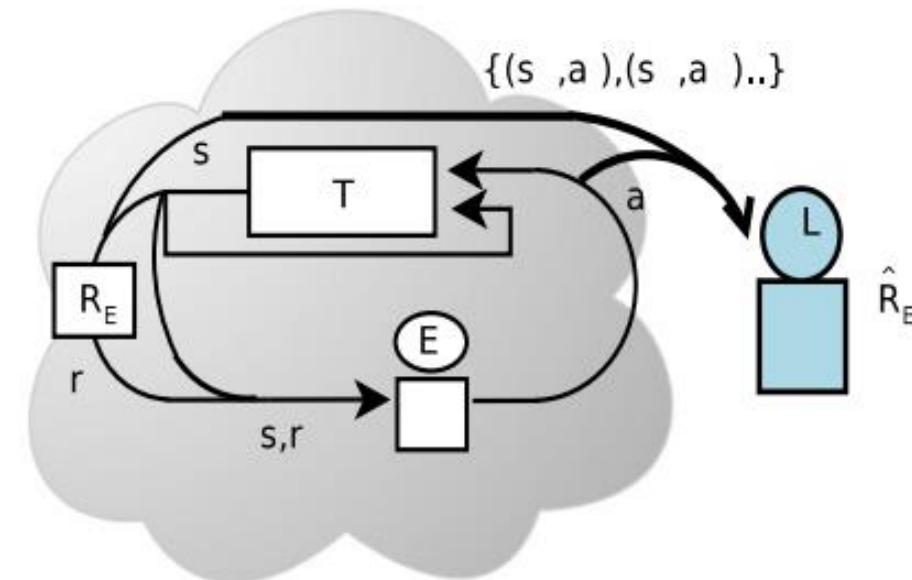
Andreas Krause
Department of Computer Science
ETH Zurich
krausea@ethz.ch

Giorgia Ramponi
ETH AI Center
giorgia.ramponi@ai.ethz.ch

IRL recovers the **feasible reward function** in which the **optimal policy** is the **expert policy**.



RL paradigm



IRL paradigm

MDP: [MDP] An MDP $\mathcal{M} := \langle S, A, T, R, \gamma \rangle$

Transition: $T : S \times A \rightarrow \text{Prob}(S)$

Reward: $R : S \times A \rightarrow \mathbb{R}$

Value: $V^\pi(s_0) = E_{s,\pi(s)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) | s_0 \right]$

MDP goal: $V^{\pi^*}(s) = V^*(s) = \sup_{\pi} V^\pi(s)$

Visit frequency: $\psi^\pi(s) = \psi^0(s) + \gamma \sum_{s' \in S} T(s, \pi(s), s') \psi^\pi(s')$

Value: $V^*(s) = \sup_{\pi} \sum_{s \in S} \psi_*^\pi(s) R(s, \pi(s))$

MDP without reward: $\mathcal{M} \setminus R$

$$\begin{aligned}\text{Feature reward: } R(s, a) &= w_1\phi_1(s, a) + w_2\phi_2(s, a) + \dots + w_k\phi_k(s, a) \\ &= \mathbf{w}^T \boldsymbol{\phi}(s, a).\end{aligned}$$

$$\text{Feature reward count: } \mu^{\phi_k}(\pi) = \sum_{t=0}^{\infty} \psi^{\pi}(s_t) \phi_k(s_t, \pi(s_t))$$

$$\text{Empirical feature count: } \hat{\mu}^{\phi_k}(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \gamma^t \phi_k(s_t, a_t).$$

$$\begin{aligned}\text{Value: } V^{\pi} &= \mathbf{w}^T \boldsymbol{\mu}^{\phi}(\pi) = \sum_{s,a} \psi^{\pi}(s) \mathbf{w}^T \boldsymbol{\phi}(s, a) \\ &= \sum_{s,a} \psi^{\pi}(s) R(s, a).\end{aligned}$$

Margin Optimization:

$$\sum_{s \in S} Q^\pi(s, a^*) - \max_{a \in A \setminus \{a^*\}} Q^\pi(s, a)$$

$$|\mu^\phi(\pi) - \hat{\mu}^\phi(\mathcal{D})|.$$

$$\hat{\pi}_E(a|s) - \pi_E(a|s)$$

Entropy Optimization:

$$\max_{\Delta} - \sum_{\tau \in (S \times A)^l} Pr(\tau) \log Pr(\tau)$$

$$\min_{P \in \Delta} \sum_{\tau \in (S \times A)^l} P(\tau) \log \frac{P(\tau)}{Q(\tau)}.$$

Bayesian Update:

$$P(\langle s, a \rangle | \hat{R}_E) \propto e^{\left(\frac{Q^*(s, a; \hat{R}_E)}{\beta} \right)}$$

Current Methods

Method	\hat{R}_E params	Optimization objective	Notable aspect
Max margin methods - maximize the margin between value of observed behavior and the hypothesis			
MMP	w	value of obs. τ - max of values from all other τ (Eq. 8)	provable convergence
MAX-MARGIN		feature exp. of policy - empirical feature exp. (Eq. 9)	sample bounds
MWAL		min diff. in value of policy and observed τ across features	first bound on iteration complexity
HYBRID-IRL		empirical stochastic policy - computed policy of expert (Eq. 10)	natural gradients and efficient optimization
LEARCH	$R(\phi)$	value of obs. τ - max of values from all other τ (Eq. 8)	nonlinear reward with suboptimal input
Silver et al. [26]			normalization of outlier inputs
Max entropy methods - maximize the entropy of the distribution over behaviors			
MAXENTIRL	w	entropy of distribution over trajectories (Eq. 11)	low learning bias
STRUCTURED APPRENTICESHIP		entropy of distribution over policies (Eq. 12)	efficient optimization
DEEP MAXENTIRL		gradient of likelihood equivalent of MaxEnt (Eq. 13)	nonlinear reward
PI-IRL			continuous state-action spaces
REIRL		relative entropy of distribution from baseline policy (Eq. 14)	suboptimal input and unknown dynamics
Bayesian learning methods - learn posterior over hypothesis space using Bayes rule			
BIRL	$R(s)$	posterior with Boltzmann data likelihood (Eq. 16)	first Bayesian IRL formulation
Lopes et al. [38]		entropy of multinomial($p_1(s), p_2(s), \dots, p_{ A -1}(s)$) derived from posterior	active learning
GP-IRL	$f(r, \theta)$	Gaussian process posterior	nonlinear reward
MLIRL	w	differentiable likelihood with Boltzmann policy (Eq. 16)	first ML approach
Classification and regression - learn a prediction model that imitates observed behavior			
SCIRL	w	Q-function as classifier scoring function	actions as state labels
CSI			provable convergence
FIRL	regression tree	norm of (\hat{R}_E - projection of \hat{R}_E)	avoids manual feature engineering

MDP\mathcal{R}

$$\mathcal{M} := (\mathcal{S}, \mathcal{A}, P, H, s_0)$$

R(s, a, h)

$$\mathcal{S} \times \mathcal{A} \times [H] \rightarrow [0, R_{\max}]$$

Q(s, a)

$$Q_{\mathcal{M} \cup r}^{\pi, h}(s, a) = r_h(s, a) + \sum_{s', a'} \pi_{h+1}(a'|s') P(s'|s, a) Q_{\mathcal{M} \cup r}^{\pi, h+1}(s', a')$$

V(s)

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

A(s, a)

$$A_{\mathcal{M} \cup r}^{\pi, h}(s, a) = Q_{\mathcal{M} \cup r}^{\pi, h}(s, a) - V_{\mathcal{M} \cup r}^{\pi, h}(s)$$

Feasible policy

$$\Pi_{\mathcal{M} \cup r}^*$$

Frequency

$$\eta_{\mathcal{M}, \pi}^{h, h}(s'|s) := \mathbb{1}_{\{s' = s\}} \text{ and } \eta_{\mathcal{M}, \pi}^{h, h'+1}(s'|s) := \sum_{s'', \tilde{a}} P(s''|s'', \tilde{a}) \pi_{h'}(\tilde{a}|s'') \eta_{\mathcal{M}, \pi}^{h, h'}(s''|s).$$

Definition: We can define a policy is optimal iff the A≤0 for any states and actions.

Definition 1 (Feasible Reward Set). A reward function r is feasible for an IRL problem (\mathcal{M}, π^E) , if and only if the expert policy π^E is optimal in $\mathcal{M} \cup r$. We call the set of all feasible reward functions $\mathcal{R}_{\mathcal{M} \cup \pi^E}$ the feasible reward set. If we estimate the transition model and expert policy from samples, we refer to the recovered feasible set $\hat{\mathcal{R}}_{\hat{\mathcal{B}}} = \mathcal{R}_{\hat{\mathcal{M}} \cup \hat{\pi}^E}$ in contrast to the exact feasible set $\mathcal{R}_{\mathcal{B}} = \mathcal{R}_{\mathcal{M} \cup \pi^E}$.

Definition 2 (Optimality Criterion). Let $\mathcal{R}_{\mathcal{B}}$ be the exact feasible set and $\hat{\mathcal{R}}_{\hat{\mathcal{B}}}$ be the feasible set recovered after observing $n \geq 0$ samples collected from \mathcal{M} and π^E . We say that an algorithm for Active IRL is (ϵ, δ, n) -correct if after n iterations with probability at least $1 - \delta$ it holds that:

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

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

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

L.

1 Def.: IRL(\mathcal{M}, π^E), \mathcal{M} 为MDP/R, π^E 为专家策略.
 ① Feasible Reward Set. $\vdash r$ is feasible $\Leftrightarrow \pi^E$ is optimal for $\mathcal{M} \cup r$.
 考虑转移和策略误差 (MDP误差, trajectories 与 π^E 的误差)
 $R_B = \mathcal{R}_{\mathcal{M} \cup \pi^E} \nearrow R_{\hat{\mathcal{B}}} = \mathcal{R}_{\hat{\mathcal{M}} \cup \hat{\pi}^E} \rightarrow$ 真的.

② 最优约束:
 $\inf_{r \in \mathcal{R}_{\mathcal{B}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup r}^*} \max_{s,a,h} |Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)| \leq \epsilon \quad \text{for each } r \in \mathcal{R}_{\mathcal{B}},$
 $\inf_{\hat{r} \in \hat{\mathcal{R}}_{\hat{\mathcal{B}}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup \hat{r}}^*} \max_{s,a,h} |Q_{\mathcal{M} \cup \hat{r}}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a)| \leq \epsilon \quad \text{for each } \hat{r} \in \hat{\mathcal{R}}_{\hat{\mathcal{B}}}.$
 $\exists r \in \mathcal{R}_{\mathcal{B}}, \forall \pi^* \in \Pi_{\mathcal{M} \cup r}^* \text{ 对于 } \forall r \in \mathcal{R}_{\mathcal{B}} \text{ 均有 } |Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup r}^{\hat{\pi}^*, h}(s, a)| < \epsilon.$
 $\exists \hat{r} \in \hat{\mathcal{R}}_{\hat{\mathcal{B}}}, \forall \pi^* \in \Pi_{\mathcal{M} \cup \hat{r}}^* \text{ 对于 } \forall \hat{r} \in \hat{\mathcal{R}}_{\hat{\mathcal{B}}} \text{ 均有 } |Q_{\mathcal{M} \cup \hat{r}}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a)| < \epsilon.$
 (recall)
 防止不必要地扩大范围. precision -

2. 有限 horizon 的约束.

Recall: 存在一个 r_{hat} 产生的策略能和 R_B 里均有相近的表现，说明查全率高。

Precision: 存在一个 R_B 里的 r 能够和 R_B_{hat} 中的所有 r_{hat} 都有相近的表现，存在限定了范围

Lemma 3 (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 4 (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:

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

4.2. 有限 horizon 的约束.

r is feasible $\Leftrightarrow A_{\mathcal{M} \cup r}^{\pi, h}(s, a) = 0$ if $\pi_h^E(a|s) \geq 0$ 且 $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) \leq 0$ if $\pi_h^E(a|s) < 0$.

Implicit: $A_{\mathcal{M} \cup r}^{\pi, h}(s, a) = r_h(s, a) + \sum_{s', a'} \pi_{h+1}(a'|s') P(s'|s, a) Q_{\mathcal{M} \cup r}^{\pi, h+1}(s', a') - \sum_a \pi_h(a|s) Q_{\mathcal{M} \cup r}^{\pi, h}(s, a)$

 $= r_h(s, a) + \sum_{s', a'} P(s'|s, a) V_{\mathcal{M} \cup r}^{\pi, h+1}(s') - V_{\mathcal{M} \cup r}^{\pi, h}(s) \rightarrow \text{Sarsa?}$

Explicit: r is feasible $\Leftrightarrow \exists \{A_h \in \mathbb{R}_{\geq 0}^{S \times A}\}_{h \in [H]} \text{ 且 } \{V_h \in \mathbb{R}^S\}_{h \in [H]}$ 对 $\forall s, a, h$ 存在.

$r_h(s, a) = \underbrace{-A_h(s, a) \mathbf{1}_{\{\pi_h^E(a|s)=0\}}}_{A_h(s, a) < 0 \text{ 且 } \pi_h^E(a|s) > 0} + V_h(s) + \sum_{s'} P(s'|s, a) V_{h+1}(s')$

不明白.

选择了但优劣较小, 则舍
对当前选择的为0, 则舍没选择的 $A_h(s, a)$

Theorem 5 (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 $A_h \leq (H - h)R_{\max}$.

专家策略误差

转移模型误差

Lemma 6. Let \mathcal{M} be an MDP \ R, r, \widehat{r} two reward functions with optimal policies $\pi^*, \widehat{\pi}^*$. Then,

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

可用Theorem5代入获得专家策略和转移模型各自的误差

$$\hat{P}_k(s'|s, a) = \frac{\sum_{h=1}^H n_k^h(s, a, s')}{\max(1, \sum_{h=1}^H n_k^h(s, a))} \quad \hat{\pi}_{k,h}^E(a|s) = \frac{n_k^h(s, a)}{\max(1, n_k^h(s))}.$$

In Appendix B.3 we derive Hoeffding's confidence intervals for the transition model and the expert policy. Combining these with Theorem 5, we can compute the uncertainty on the recovered reward as:

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

where $\ell_k^h(s, a) = \log(24SAH(n_k^h(s, a))^2/\delta)$. We can show that for any pair of reward functions $r \in \mathcal{R}_{\mathfrak{B}}$ and $\hat{r} \in \mathcal{R}_{\hat{\mathfrak{B}}}$, the difference $|r_h(s, a) - \hat{r}_{k,h}(s, a)| \leq C_k^h(s, a)$. This uncertainty estimate will be a key component in all of our theoretical analysis.

Theorem 7 (Sample Complexity of Uniform Sampling IRL). *The uniform sampling strategy fulfills Definition 2 with a number of samples upper bounded by:*

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

where \mathcal{O} suppresses logarithmic terms.

Ace-IRL without generative model

Definition 2 (Optimality Criterion). 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 an algorithm for Active IRL is (ϵ, δ, n) -correct if after n iterations with probability at least $1 - \delta$ it holds that:

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

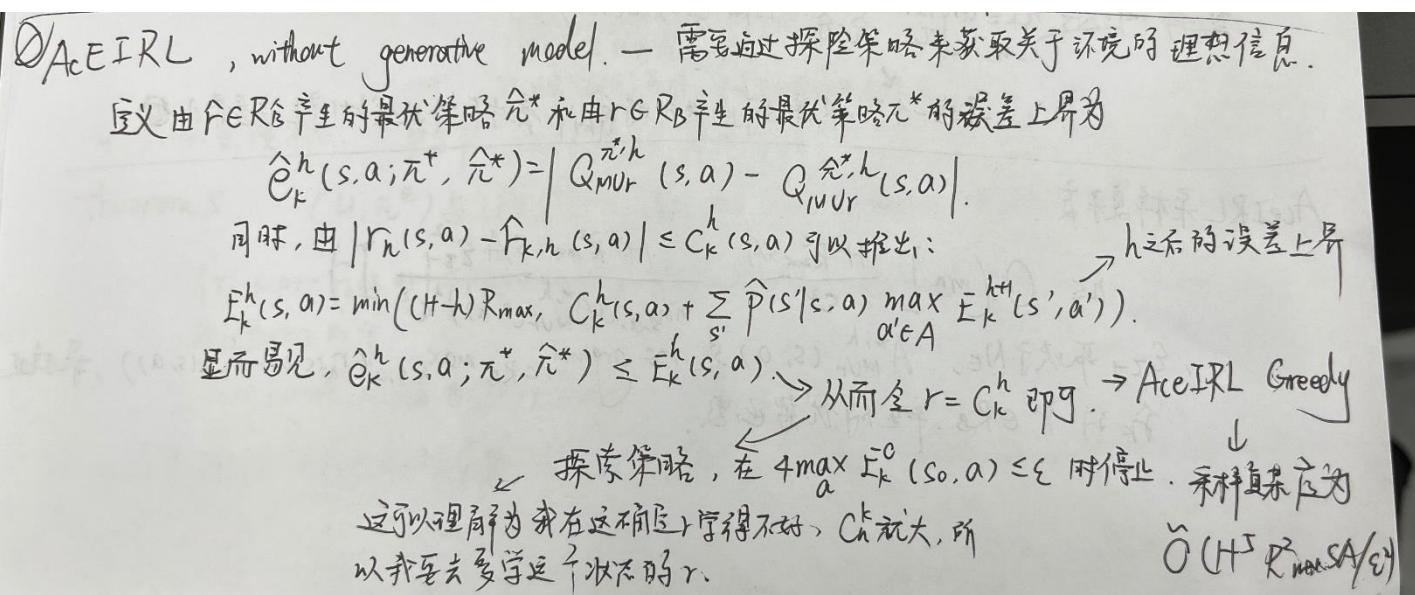
$$\inf_{r \in \mathcal{R}_{\mathfrak{B}}} \sup_{\pi^* \in \Pi_{\mathcal{M} \cup r}^*} \max_{s, a, h} \left| Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) \right| \leq \epsilon \quad \text{for each } \hat{r} \in \mathcal{R}_{\hat{\mathfrak{B}}},$$

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

$$E_k^h(s, a) = \min \left((H - h) R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{P}(s' | s, a) \max_{a' \in \mathcal{A}} E_k^{h+1}(s', a') \right). \quad (\text{EB1}) \quad 4 \max_a E_k^0(s_0, a) \leq \epsilon.$$

$$\hat{e}_k^h(s, a; \pi^*, \hat{\pi}^*) = \left| Q_{\mathcal{M} \cup r}^{\pi^*, h}(s, a) - Q_{\mathcal{M} \cup \hat{r}}^{\hat{\pi}^*, h}(s, a) \right|.$$

Error Setting



在这里已经可以单纯的将误差上界作为奖赏去进行探索了, 可以理解为误差上界越高代表当前的状态信息越少, 越需要更多的信息多 reward 进行更好的近似。

这种探索策略称为 AceIRL-Greedy。

Drawbacks of the AceIRL-greedy:

1. AceIRL-greedy doesn't reduce the future uncertainty.

AceIRL 缺陷：① 探索了高不确定性状态，但我们的目标在于在下一次迭代减少不确定性。
② AceIRL 探索为减少所有策略的不确定性，但目标是减少合理的策略和其他策略的不确定性。

所以作者提出我们需要一个在探索过程中减少不确定性的探索策略，即我们需要一个能够最大化 E_{k+1}^h 的探索策略。但 E_{k+1}^h 不好算，可以避对当前转移模型的估计。

① $\hat{C}_{k+1}^h(s, a) = (H-h) R_{\max} \min(1, 2 \sqrt{\frac{\hat{\pi}_k^h(s, a)}{\hat{\pi}_k^h(s, a) + \hat{\pi}_{\pi_k^h}(s, a)}})$.

$\hat{\pi}_{\pi_k^h}(s, a) = N_E \cdot \eta_{M, \pi_k^h}^{0, h}(s, a | s_0) \rightarrow \pi$ 在时间 t 访问 (s, a) 的期望数。
× N_E 为使用元探索的次数。

尽管我们的目标是不确定性采样，但现在这个能更好的衡量选择一个探索策略的信息是这个问题不基于 IRL，并可用于提升 reward-free 探索。

Drawbacks of the AceIRL-greedy:

2. AceIRL-greedy reduce all policies' uncertainty.

AceIRL 缺陷：① 探索了高不确定性状态，但我们的目标在于在下一次迭代减少不确定性。

② AceIRL 探索为减少所有策略的不确定性，但目标是减少合理的其他策略的不确定性。

我们只对 $\pi^* \in \Pi_{\text{MUR}}$ 和 $\hat{\pi}^* \in \Pi_{\text{MUR}}^*$ 做更新， $r \in R_B$, $\hat{r} \in R_B^*$.

假设我们能构造出一批可靠最优策略 $\hat{\pi}_k$ ，其中包含有 π^* , $\hat{\pi}_k^*$ (大概率)。

重新取上界 \nearrow h 步后的误差上界

$$\hat{\pi}_k^h(s, a) = 0$$

$$\hat{\pi}_k^h(s, a) = \min((H-h)R_{\max}, C_k^h(s, a) + \sum_{s'} \hat{\phi}(s'|s, a) \max_{\pi \in \Pi_{k+1}} \pi(a'|s') \hat{\pi}_{k+1}^{h+1}(s', a')).$$

与前面的区别在于 $\hat{\pi}_k^h$ 中的乘积而非所有动作。

$$\text{停止条件 } \max_a \hat{\pi}_k^0(s_0, a) \leq \varepsilon. \rightarrow R_B \text{ 满足 Def 2.}$$

Theorem 8. [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 .

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. This bound matches the bound we saw previously with a generative model. Hence, AceIRL achieves the same results without access to the generative model. 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. If N_E is small, then the ϵ -ratio will be also small. If N_E is large the algorithm will perform similarly to a uniform sampling strategy. Appendix B.5 provides the full proof of this theorem.

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

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

Algorithm 1 AceIRL algorithm for IRL in an unknown environment.

- 1: **Input:** significance $\delta \in (0, 1)$, target 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: $k \leftarrow k + 1$
 - 7: Update \hat{P}_k , $\hat{\pi}_k$, C_k^h , and $\hat{r}_k \leftarrow \mathcal{A}(\mathcal{R}_{\hat{\mathcal{B}}})$
 - 8: Update accuracy $\epsilon_k \leftarrow \max_a \hat{E}_k^0(s_0, a)$
 - 9: **end while**
 - 10: **return** Estimated reward function \hat{r}_k
-

Experiment

	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 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 C.3, we show learning curves for all individual experiments.

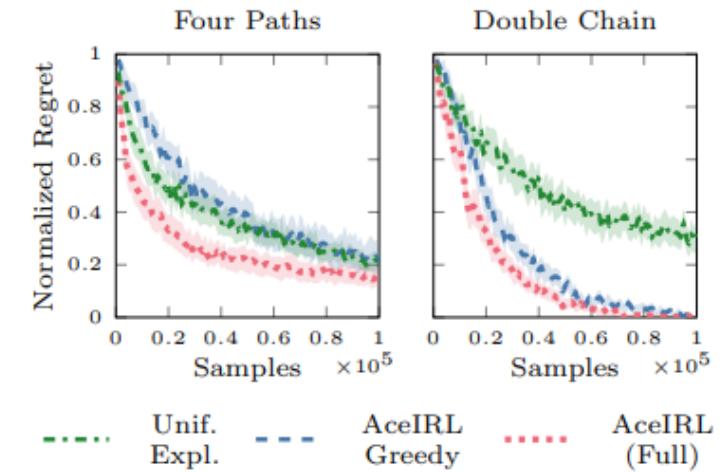


Figure 2: Normalized regret (lower is better) of the policy optimizing for the inferred reward in the estimated MDP as a function of the number of samples. The plots show the mean and 95% confidence intervals computed using 50 random seeds. We use $N_E = 50$.

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 $\mathcal{M} \cup (-r)$.

Thanks