Efficient Algorithms for Inverse Reinforcement Learning

Gokul Swamy

Recap from Last Time

- 1. Why do we need interaction in imitation learning?
- **A**: to be able to tell that we've made a mistake that compounds.
- 2. What else do we need to tell which mistakes matter?
- A: information about the set of rewards we could be judged on.
- 3. How do we learn a policy that recovers from mistakes that matter if we don't know what the reward function is?
- **A**: Find the policy that is the least distinguishable from the expert's under any reward function in the moment set \mathcal{R} .



1. What makes inverse RL sample-inefficient?

2. Are best responses required for solving the IRL game?

3. What algorithms can we use in our new reduction?

1. What makes inverse RL sample-inefficient?

A: Repeatedly solving a global exploration (RL) problem.

3. What algorithms can we use in our new reduction?

- 2. Are best responses required for solving the IRL game?

Recap: Inverse RL as Game-Solving

$\max_{\pi \in \Pi} \min_{f \in \mathscr{R}} J(\pi, f) - J(\pi_E, f)$

where $J(\pi, f) \triangleq \mathbb{E}_{\xi \sim \pi} \left| \sum_{k \in \pi} f(s_h, a_h) \right|.$ h



Recap: Two Flavors of IRL Algorithms $\max_{\pi \in \Pi} \min_{f \in \mathscr{R}} J(\pi, f) - J(\pi_E, f)$ Dual Primal Policy Update BR: RL NR: GD NR: 60 NR: GD **Reward Update**







We've reduced the "easier" problem of IL to the "harder" problem of RL



$\begin{array}{l} \max & \min \\ \pi \in \Pi & f \in \mathscr{R} \end{array}$	JCT
	Dı
Policy Update	BR:
Reward Update	NR

Primal algorithms also need to explore as we can force any no-regret algorithm to compute π^* by playing the same f!



- 1. What makes inverse RL sample-inefficient?
- A: Repeatedly solving a global exploration (RL) problem.
- 2. Are best responses required for solving the IRL game?
- A: Actually, competing with the expert is "all you need".
- 3. What algorithms can we use in our new reduction?



... which means it can't be the policy we want!

π_F need not be optimal for f_t ... f_2

0

+1





\sim Idea: save compute by competing with π_E , not π_t^{\star} !



Reducing IRL to Expert-Competitive RL

sequence of policies $\pi_{t+1} = A_{\pi}(f_{1,t})$ such that t = 1

 \sim Idea: We never need to compute a best response to an $f_{f}!$

ERROr { $\operatorname{Reg}_{\pi}(T)$ }: A policy-selection algorithm \mathbb{A}_{π} satisfies the $\operatorname{Reg}_{\pi}(T)$ expert-relative regret guarantee if given any sequence of reward functions $f_{1,T}$, it produces a

$\sum J(\pi_E, f_t) - J(\pi_t, f_t) \leq \operatorname{Reg}_{\pi}(T).$

Reducing IRL to Expert-Competitive RL A_f is a **no-regret reward selection algorithm** if when given a sequence of policies $\pi_{1:t}$, it produces iterates $f_{t+1} = A_f(\pi_{1:t})$ such that $\sum_{t=1}^{T} \overline{J(\pi_t, f_t)} - J(\pi_E, f_t) - \min_{f^* \in \mathcal{F}_r} \sum_{t=1}^{T} J(\pi_t, f^*) - J(\pi_E, f^*) \le H \operatorname{Reg}_f(T),$ with $\lim_{t \to \infty} \frac{\text{Reg}_f(T)}{T} = 0.$ $T \rightarrow \infty$

Reducing IRL to Expert-Competitive RL $J(\pi_E, r) - J(\bar{\pi}, r) = \frac{1}{T} \sum_{t=1}^{T} J(\pi_E, r) - J(\pi_t, r)$ reward realizability $\leq \max_{f^{\star} \in \mathcal{F}_{r}} \frac{1}{T} \sum_{t=1}^{I} J(\pi_{E}, f^{\star}) - J(\pi_{t}, f^{\star}) \quad \operatorname{clefn.of}_{\text{(equal)}}$ $\leq \frac{1}{T} \sum_{t=1}^{T} J(\pi_{E}, f_{t}) - J(\pi_{t}, f_{t}) + \frac{\operatorname{Reg}_{f}(T)}{T} H$ $\operatorname{Reg}_{f}(I)$ $\operatorname{Reg}(I)$ $\mathcal{O}_{\pi^{\mathsf{T}}}$

- 1. What makes inverse RL sample-inefficient?
- A: Repeatedly solving a global exploration (RL) problem.
- 2. Are best responses required for solving the IRL game?
- A: Actually, competing with the expert is "all you need".
- 3. What algorithms can we use in our new reduction?
- A: A wide variety of sample-efficient "local search" algorithms.

Expert-Competitive RL via PSDP



Idea: Reset to states from the demonstrations!

Poly (H, log ITI) VC-dim.





Expert-Competitive RL via PSDP

- $\pi_{H} = \arg \max_{\pi \in \Pi} \mathbb{E}_{s_{H} \sim \pi_{E}} [\mathbb{E}_{a_{H} \sim \pi} [f_{t}(s_{H}, a_{H})]]$
- For $h \in [H 1, H 2, ..., 1]$:

 $\pi_{h} = \arg\max_{\pi \in \Pi} \mathbb{E}_{s_{h} \sim \pi_{E}} [\mathbb{E}_{a_{h} \sim \pi} [f_{t}(s_{h}, a_{h}) + \mathbb{E}_{T(s_{h}, a_{h})} [V_{t}^{\pi_{h+1:H}}(s_{h+1})]]]$

Lemma: Assume that at each time-step $h \in [H]$, we perform policy optimization up to ε -optimality: $\mathbb{E}_{S_{h},a_{h}\sim\pi_{F}}[Q_{t}^{\pi_{h+1}:H}(S_{h},a_{h}) - \mathbb{E}_{a\sim\pi_{h}(S_{h})}[Q_{t}^{\pi_{h+1}:H}(S_{h},a)]] \leq \varepsilon H$ Then,

 $J(\pi_E, f_t) - J(\pi_{1:H}, f_t) \leq \mathcal{O}(\varepsilon H^2)$

Proof: We proceed via the PDL (shocking, I know):

 $J(\pi_E, f_t) - J(\pi_{1:H}, f_t) = \sum_{s_h, a_h \sim \pi_E} \left[Q^{\pi_{h+1:H}}(s_h, a_h) - \mathbb{E}_{a \sim \pi_h(s_h)}[Q^{\pi_{h+1:H}}(s_h, a_h)] \right]$ $\leq \sum_{k=1}^{n} \varepsilon(H-h)$ $\leq \mathcal{O}(\varepsilon H^2)$ $\Rightarrow \operatorname{Reg}_{\pi}(T) \leq \varepsilon H^2 T$

Unavoidable in general on cliff-like (irrecoverable) problems.





Interaction can still help us figure out which mistakes compound. (Can interpolate with or anneal towards ρ_0 in practice)





poly(H) Algorithms for IRL! -> Re wrosot. / to extert states ~, SAC **E** π_{t+2} Idea: Localize policy search by resetting to states from the demonstrations!









- 1. What makes inverse RL sample-inefficient?
- A: Repeatedly solving a global exploration (RL) problem.
- 2. Are best responses required for solving the IRL game?
- A: Actually, competing with the expert is "all you need".
- 3. What algorithms can we use in our new reduction?
- A: A wide variety of sample-efficient "local search" algorithms. Lo NRPI, PSOP, Hya, Agroste Storf)



- $\pi_{H} = \arg\max_{\pi \in \Pi} \mathbb{E}_{s_{H} \sim \mu} [\mathbb{E}_{a_{H} \sim \pi} [f_{t}(s_{H}, a_{H})]]$
- For $h \in [H 1, H 2, ..., 1]$:
 - $\pi_h = \arg \max_{\pi \in \Pi} \mathbb{E}_{s_h \sim \mu} [\mathbb{E}_{a_h \sim \pi} [f]]$
- policy optimization up to ε -optimality: I hen,

 $J(\pi'_{1:H}, f_t) - J(\pi_{1:H}, f_t) \le \mathcal{O}((\varepsilon + \mathbb{D}_{TV}(\mu, \rho_{\pi'})) \cdot H^2)$

If $\textcircled{O}: PSDP Competes Against Policies Covered by <math>\mu$

$$f_t(s_h, a_h) + \mathbb{E}_{T(s_h, a_h)}[V_t^{\pi_{h+1}:H}(s_{h+1})]$$

Lemma: Assume that at each time-step $h \in [H]$, we perform

 $\mathbb{E}_{s_{h} \sim \mu}[\mathbb{E}_{a \sim \pi'_{h}(s_{h})}[Q_{t}^{\pi_{h+1}:H}(s_{h},a)] - \mathbb{E}_{a \sim \pi_{h}(s_{h})}[Q_{t}^{\pi_{h+1}:H}(s_{h},a)]] \leq \varepsilon H$

Proof: We proceed via the PDL:

 $J(\pi'_{1:H}, f_t) - J(\pi_{1:H}, f_t) = \sum_{h=1}^{H} \mathbb{E}_{s_h, a_h \sim \pi'_{1:h}} \left[Q^{\pi_{h+1:H}}(s_h, a_h) - \mathbb{E}_{a \sim \pi_h(s_h)}[Q^{\pi_{h+1:H}}(s_h, a_h)] \right]$ $\leq \sum_{k=1}^{n} \mathbb{E}_{s_{h},a_{h}\sim\mu} \left[Q^{\pi_{h+1:H}}(s_{h},a_{h}) - \mathbb{E}_{a\sim\pi_{h}(s_{h})}[Q^{\pi_{h+1:H}}(s_{h},a_{h})] \right]$ + $H \cdot \mathbb{D}_{TV}(\mu_h, \rho_h^{\pi'})$ $\leq \sum \left(\varepsilon + \mathbb{D}_{TV}(\mu_h, \rho_h^{\pi'})\right) \cdot 2 \cdot (H - h)$ $\leq (\varepsilon + \mathbb{D}_{TV}(\mu, \rho_{\pi'})) \cdot H^2$





If 3: What μ should we reset to when $\pi_E \notin \Pi$?



Widen the baseline distribution μ to cover $\pi^* \in \Pi$, potentially by using suboptimal / offline data!