



POLITECNICO
MILANO 1863

Provably Efficient Learning of Transferable Rewards

Alberto Maria Metelli* Giorgia Ramponi* Alessandro Concetti Marcello Restelli

July 2021

Thirty-eighth International Conference on Machine Learning

- **Goal:** learn *one* reward function from **expert's** demonstrations (Ng and Russell, 2000)



- IRL problem is **ambiguous** (Abbeel and Ng, 2004)!
- **Feasible Reward Set:** set of all rewards making π^E optimal (Ng and Russell, 2000)

$$\mathcal{R} = \{r \in \mathbb{R}^{S \times A} : \pi^E \in \text{Greedy}(Q_{M \cup r}^*)\}$$

- **Goal:** learn *one* reward function from **expert's** demonstrations (Ng and Russell, 2000)



- IRL problem is **ambiguous** (Abbeel and Ng, 2004)!
- **Feasible Reward Set:** set of all rewards making π^E optimal (Ng and Russell, 2000)

$$\mathcal{R} = \{r \in \mathbb{R}^{S \times A} : \pi^E \in \text{Greedy}(Q_{M \cup r}^*)\}$$

- **Goal:** learn *one* **reward** function from **expert**'s demonstrations (Ng and Russell, 2000)



- IRL problem is **ambiguous** (Abbeel and Ng, 2004)!
- **Feasible Reward Set:** set of all rewards making π^E optimal (Ng and Russell, 2000)

$$\mathcal{R} = \{r \in \mathbb{R}^{S \times A} : \pi^E \in \text{Greedy}(Q_{M \cup r}^*)\}$$

- **Goal:** learn *one* reward function from **expert's** demonstrations (Ng and Russell, 2000)



- IRL problem is **ambiguous** (Abbeel and Ng, 2004)!
- **Feasible Reward Set:** set of all rewards making π^E optimal (Ng and Russell, 2000)

$$\mathcal{R} = \{r \in \mathbb{R}^{S \times A} : \pi^E \in \text{Greedy}(Q_{M \cup r}^*)\}$$

- (Q1) How does the error on the **transition model** \hat{P} and on the **expert's policy** $\hat{\pi}^E$ propagate to the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$?
- (Q2) How does the error on the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$ affect the **value function** $Q_{\mathcal{M}' \cup \hat{r}}^*$ in a **different environment** \mathcal{M}' ?



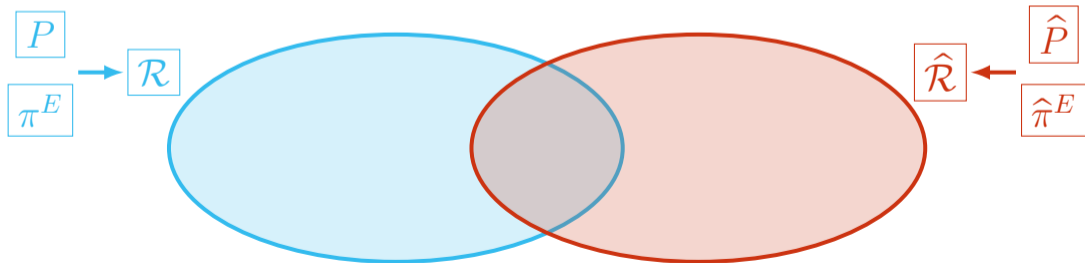
- (Q1) How does the error on the **transition model** \hat{P} and on the **expert's policy** $\hat{\pi}^E$ propagate to the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$?
- (Q2) How does the error on the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$ affect the **value function** $Q_{\mathcal{M}' \cup \hat{r}}^*$ in a **different environment** \mathcal{M}' ?



(Q1) Reward Error Propagation

(Q1) How does the error on the **transition model** \hat{P} and on the **expert's policy** $\hat{\pi}^E$ propagate to the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$?

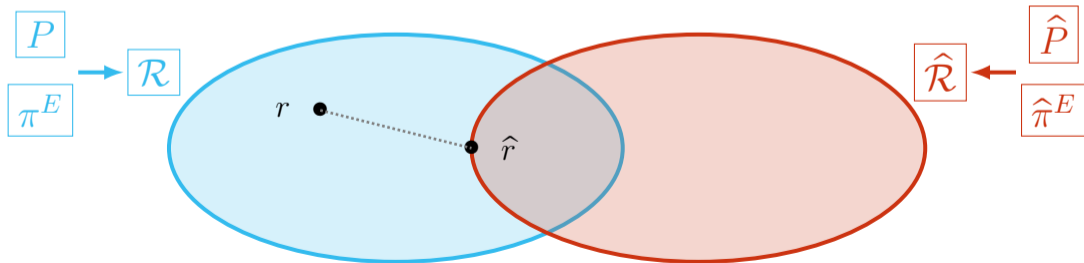
$$\inf_{\hat{r} \in \hat{\mathcal{R}}} |r - \hat{r}| \leq \frac{R_{\max}}{1 - \gamma} \text{Distance}(\pi^E, \hat{\pi}^E) + \frac{\gamma R_{\max}}{1 - \gamma} \text{Distance}(P, \hat{P})$$



(Q1) Reward Error Propagation

(Q1) How does the error on the **transition model** \hat{P} and on the **expert's policy** $\hat{\pi}^E$ propagate to the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$?

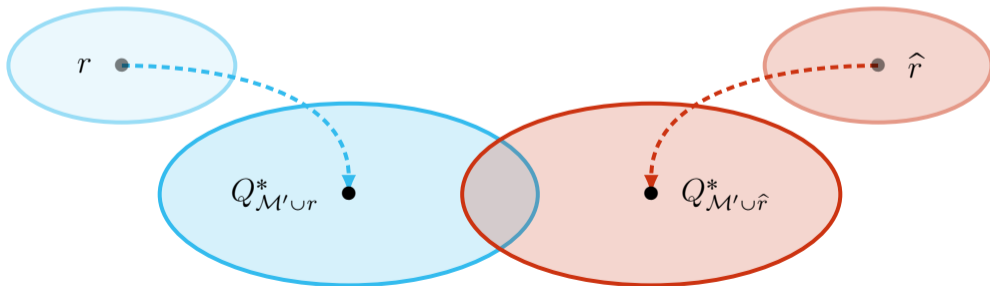
$$\inf_{\hat{r} \in \hat{\mathcal{R}}} |r - \hat{r}| \leq \frac{R_{\max}}{1 - \gamma} \text{Distance}(\pi^E, \hat{\pi}^E) + \frac{\gamma R_{\max}}{1 - \gamma} \text{Distance}(P, \hat{P})$$



(Q2) Transferred Reward Error Propagation

(Q2) How does the error on the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$ affect the **value function** $Q_{\mathcal{M}' \cup \hat{r}}^*$ in a **different environment** \mathcal{M}' ?

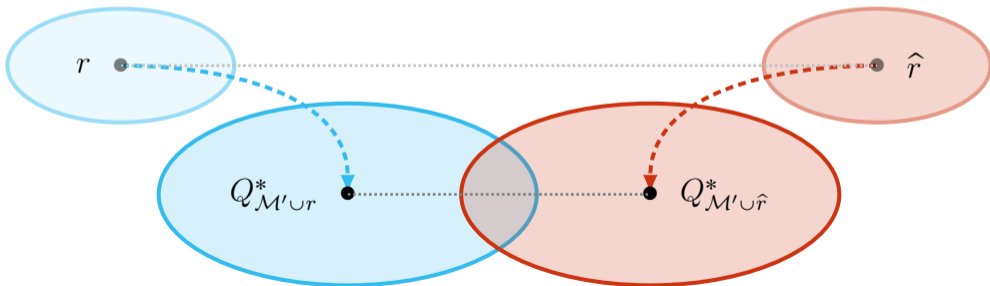
$$\|Q_{\mathcal{M}' \cup r}^* - Q_{\mathcal{M}' \cup \hat{r}}^*\|_{\infty} \leq \max_{\pi \in \{\pi^*, \hat{\pi}^*\}} \underbrace{\| (I_{S \times A} - \gamma' P' \pi)^{-1} \|}_{\text{occupancy in } \mathcal{M}'}} (r - \hat{r})\|_{\infty}$$



(Q2) Transferred Reward Error Propagation

(Q2) How does the error on the **recovered reward** $\hat{r} \in \hat{\mathcal{R}}$ affect the **value function** $Q_{\mathcal{M}' \cup \hat{r}}^*$ in a **different environment** \mathcal{M}' ?

$$\|Q_{\mathcal{M}' \cup r}^* - Q_{\mathcal{M}' \cup \hat{r}}^*\|_{\infty} \leq \max_{\pi \in \{\pi^*, \hat{\pi}^*\}} \underbrace{\| (I_{S \times A} - \gamma' P' \pi)^{-1} (r - \hat{r}) \|_{\infty}}_{\text{occupancy in } \mathcal{M}'}$$



- Problem Setting

- **Generative model** of $\mathcal{M} = (P, \gamma)$ and possibility of sampling from π^E
- Target MDP $\mathcal{M}' = (P', \gamma')$ **known**

- (ϵ, δ, n) -correct **Sampling Strategy**

- Fix target rewards $(\bar{r}, \check{r}) \in \mathcal{R} \times \hat{\mathcal{R}}$
- After having collected n samples

$$\inf_{\hat{r} \in \hat{\mathcal{R}}} \|Q_{\mathcal{M}' \cup \bar{r}}^* - Q_{\mathcal{M}' \cup \hat{r}}^*\|_{\infty} \leq \epsilon \quad \text{and} \quad \inf_{r \in \mathcal{R}} \|Q_{\mathcal{M}' \cup r}^* - Q_{\mathcal{M}' \cup \check{r}}^*\|_{\infty} \leq \epsilon \quad \text{w.p. } 1 - \delta$$

- Problem Setting

- **Generative model** of $\mathcal{M} = (P, \gamma)$ and possibility of sampling from π^E
- Target MDP $\mathcal{M}' = (P', \gamma')$ **known**

- (ϵ, δ, n) -correct **Sampling Strategy**

- Fix target rewards $(\bar{r}, \check{r}) \in \mathcal{R} \times \hat{\mathcal{R}}$
- After having collected n samples

$$\inf_{\hat{r} \in \hat{\mathcal{R}}} \|Q_{\mathcal{M}' \cup \bar{r}}^* - Q_{\mathcal{M}' \cup \hat{r}}^*\|_{\infty} \leq \epsilon \quad \text{and} \quad \inf_{r \in \mathcal{R}} \|Q_{\mathcal{M}' \cup r}^* - Q_{\mathcal{M}' \cup \check{r}}^*\|_{\infty} \leq \epsilon \quad \text{w.p. } 1 - \delta$$

- Allocate samples **uniformly** over states and actions
- (ϵ, δ, n) -correct with $n = \sum_{s,a} n(s, a)$:

$$n(s, a) \leq \tilde{O} \left(\frac{\gamma^2 R_{\max}^2}{(1-\gamma)^2 (1-\gamma)^2 \epsilon^2} \right)$$

- Allocate samples **uniformly** over states and actions
- (ϵ, δ, n) -correct with $n = \sum_{s,a} n(s, a)$:

$$n(s, a) \leq \tilde{O} \left(\frac{\gamma^2 R_{\max}^2}{(1 - \gamma')^2 (1 - \gamma)^2 \epsilon^2} \right)$$

- Allocate samples based on the **error propagation bound**
- (ϵ, δ, n) -correct with $n = \sum_{s,a} n(s, a)$:

$$n(s, a) \leq \tilde{O} \left(\min \left\{ \frac{\gamma^2 R_{\max}^2}{(1-\gamma)^2 (1-\gamma)^2 \epsilon^2}, \frac{\gamma^2 R_{\max}^2 \sigma^2}{(1-\gamma)^2 (1-\gamma)^2 \epsilon^2} \right\} \right)$$

$$\Delta(s, a) = V_{\gamma}^*(s) - Q_{\gamma}^*(s, a)$$

- Allocate samples based on the **error propagation bound**
- (ϵ, δ, n) -correct with $n = \sum_{s,a} n(s, a)$:

$$n(s, a) \leq \tilde{O} \left(\min \left\{ \frac{\gamma^2 R_{\max}^2}{(1 - \gamma')^2 (1 - \gamma)^2 \epsilon^2}, \frac{\gamma^2 R_{\max}^2 (\epsilon')^2}{(1 - \gamma)^2 \Delta(s, a)^2 \epsilon^2} \right\} \right)$$

$$\Delta(s, a) = V_{\tilde{r}}^*(s) - Q_{\tilde{r}}^*(s, a)$$

Thank You for Your Attention!

Code: `github.com/albertometelli/travel`

Contact: `albertomaria.metelli@polimi.it`



- P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In C. E. Brodley, editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004), Banff, Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM, 2004. doi: 10.1145/1015330.1015430.
- A. Y. Ng and S. J. Russell. Algorithms for inverse reinforcement learning. In P. Langley, editor, *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pages 663–670. Morgan Kaufmann, 2000.