# Approximate Robust Inverse Reinforcement Learning

Brendan Crowe

- **Reinforcement learning (RL)** is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. [2]

- **Reinforcement learning (RL)** is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. [2]

- The optimizer knows that reward function $r$

- **Reinforcement learning (RL)** is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. [2]

- The optimizer knows that reward function $r$

- The optimizer wants to find the policy $\pi^\star \in \Pi$

- **Reinforcement learning (RL)** is an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. [2]

- The optimizer knows that reward function $r$

- The optimizer wants to find the policy $\pi^\star \in \Pi$

- Example: AlphaGO

- **Inverse Reinforcement Learning (IRL):** Does not assume that the reward function *r* is known, instead the goal is to find a reward function that best explains some expert data.

- **Inverse Reinforcement Learning (IRL):** Does not assume that the reward function *r* is known, instead the goal is to find a reward function that best explains some expert data.
- The optimizer has data $\mathcal{D}$ that shows the tasks being done (hopefully) optimally

- **Inverse Reinforcement Learning (IRL):** Does not assume that the reward function *r* is known, instead the goal is to find a reward function that best explains some expert data.
- The optimizer has data $\mathcal{D}$ that shows the tasks being done (hopefully) optimally
- The optimizer wants to find a reward function $r^* \in R$ that best explains $\mathcal{D}$ and produces a good policy $\pi^*$

- **Inverse Reinforcement Learning (IRL):** Does not assume that the reward function *r* is known, instead the goal is to find a reward function that best explains some expert data.
- The optimizer has data $\mathcal{D}$ that shows the tasks being done (hopefully) optimally
- The optimizer wants to find a reward function $r^\star \in R$ that best explains $\mathcal{D}$ and produces a good policy $\pi^\star$
- Example: Teaching a self driving car using data collected from people driving in the real world

- Infinite number of reward functions that fit the data

- The goal is to be able to find a policy that does well for the worst-case reward function

- We can deal with this uncertainty using robust optimization

$$\max_{u} \min_{r} \quad r^T u$$
$$s.\,t. \quad A^T u = c$$
$$u \geq 0 \qquad [1]$$

(1)

$$\max_{u,z} \quad z$$
$$s.t. \quad A^T u = c$$
$$z \leq r^T u \quad r \in R \qquad\qquad (2)$$
$$u \geq 0 \qquad [1]$$

- $R$: a set of possible reward functions
- $u$: the occupancy for each state-action pair
- $A = \begin{bmatrix} I - \gamma P_a \\ \vdots \end{bmatrix}$
- $P_a$: probability of transitioning from one state to the next given an action
- $0 \leq \gamma \leq 1$: discount rate (a constant)
- $c$: the distribution over starting states-actions

- This requires all states to be known

- Not feasible for large problems

- Not flexible

- The solution: approximation

- Use linear features $\Phi$ that approximate the states

- Generalizes to states we haven't seen

$$\max_{u,z} \quad z$$
$$s.\,t. \quad \Phi A^T u = \Phi c$$
$$\qquad z \leq r^T u \quad r \in R \tag{3}$$
$$\qquad u \geq 0$$

What's going on??

$$\min_{w,\xi} \quad c^T \Phi w$$
$$s.t. \quad \Phi A^T w \geq \hat{\Phi}\tilde{R}\xi \tag{4}$$
$$\mathbb{1}^T \xi = 1$$
$$\xi \geq 0$$

- $\Phi w \approx v$: represents the value of being at each state
- $\Phi \in \mathbb{R}^{n \times m}$: a features matrix that is number of observed states $\times$ number of features
- $\hat{\Phi} = \begin{bmatrix} \Phi & 0 \\ 0 & \Phi \end{bmatrix}$
- $\tilde{R}$: weights that describe inform a reward function $\hat{\Phi}\tilde{R} \approx R$
- $\xi$: a weight for each possible reward function in $\hat{\Phi}\tilde{R}$

[1] Daniel S. Brown, Scott Niekum **and** Marek Petrik. *Bayesian Robust Optimization for Imitation Learning*. 2020. arXiv: `2007.12315 [cs.LG]`.

[2] Wikipedia contributors. *Reinforcement learning — Wikipedia, The Free Encyclopedia*. [Online; accessed 3-May-2021]. 2021. URL: `https://en.wikipedia.org/w/index.php?title=Reinforcement_learning&oldid=1019990151`.