

# Lecture 21 - Imitation learning

Guiliang Liu

The Chinese University of Hong Kong, Shenzhen

DDA4230: Reinforcement Learning

Course Page: [\[Click\]](#)

# Imitation Learning

**Motivation.** Learning policies from rewards is successful in situations where data is cheap and easily gathered. This approach fails, however, when **data gathering is slow**, **failure must be avoided** (e.g. autonomous vehicles), or **safety is desired**.

- One approach to mitigate the sparse reward problem is to **manually design reward functions** that are dense in time. However, this approach requires a human to hand-design a reward function with the desired behavior in mind.
- It is therefore desirable to learn by **imitating agents performing** the task in question.



# Imitation Learning

Generally, experts provide a set of demonstration trajectories, which are sequences of states and actions. More formally, we assume that we are given

- State space, action space;
- Access to the transition oracle  $\mathbb{P}(s' | s, a)$ ;
- Set of one or more teacher demonstrations  $(s_0, a_0, s_1, a_1, \dots)$ , where actions are drawn from the teacher's policy  $\pi^*$ .

However, **no reward function** oracle  $\mathcal{R}$  and **no explicit transition** model  $\mathbb{P}(s' | s, a)$  are given.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Behavioral Cloning

A natural question raised out of this context is then

*Can we learn the teacher's policy using supervised learning?*

In behavioral cloning, we aim simply to learn the policy via supervised learning.

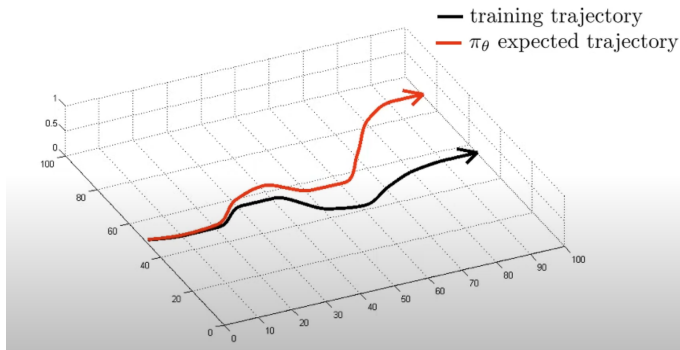
- Specifically, we will fix a policy class and aim to learn a policy mapping states to actions given the data tuples  $\{(s_0, a_0), (s_1, a_1), \dots\}$ .



# Behavioral Cloning

One challenge to this approach is that data is not distributed i.i.d. in the state space.

In RL, errors are compounding and they accumulate over the length of the episode.

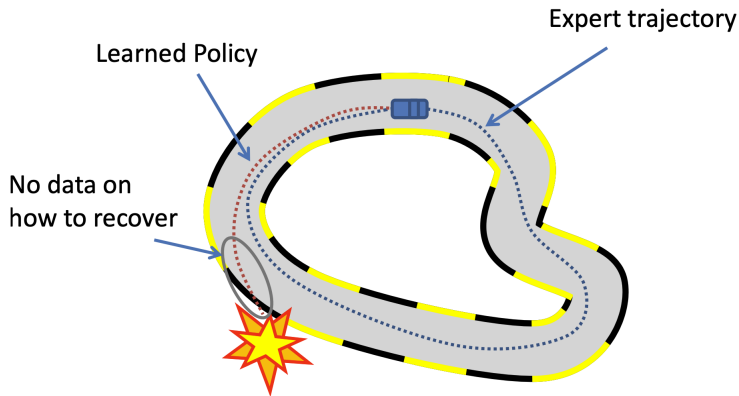


中文大學(深圳)

Chinese University of Hong Kong, Shenzhen

# Behavioral Cloning

One challenge to this approach is that data is not distributed i.i.d. in the state space.  
In RL, errors are compounding and they accumulate over the length of the episode.



# Behavioral Cloning

One challenge to this approach is that data is not distributed i.i.d. in the state space. In RL, **errors are compounding and they accumulate over the length of the episode.**

- The training data for the learned policy will be tightly clustered around expert trajectories.
- If a mistake is made that puts the agent in a part of the state space that the expert did not visit, the agent has no data to learn a policy from.
- The error scales quadratically in the episode length, as opposed to the linear scaling in standard RL.



# Behavioral Cloning

## DAGGER: Dataset aggregation:

- This algorithm aims to mitigate the problem of compounding errors by adding data for newly visited states.
- As opposed to assuming there is a pre-defined set of expert demonstrations, we assume that we can generate more data from an expert.
- The limitation of this, of course, is that an expert must be available to provide labels, sometimes in real-time.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen



# Behavioral Cloning

---

**Algorithm 1: DAGGER**

---

Initialize  $\mathcal{D} \leftarrow \emptyset$

Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$

**for**  $i = 1$  *to*  $N$  **do**

    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$

    Sample  $T$ -step trajectories using  $\pi_i$

    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$  and actions given by expert

    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$

    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$

**return** best  $\hat{\pi}_i$  on validation

---



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Inverse Reinforcement Learning

**Motivation.** Behavior cloning directly learns the policy as desired, but its practical performance can be limited. The reason is that apart from the input provided by the experts, **there are not many generalizations that are provided by the algorithm**. Instead, a better generalization can be obtained by learning the **reward function**, which is a succinct description of the task, from the expert input.

*Can we recover the reward function  $\mathcal{R}$  from expert input?*



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Inverse Reinforcement Learning

*Can we recover the reward function  $\mathcal{R}$  from expert input?*

In inverse reinforcement learning, the goal is to **learn the reward function** (that has not been provided) based on the expert demonstrations.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Inverse RL via Margin Optimization

We consider a reward which is represented as a linear combination of features

$$R(s) = w^T x(s),$$

where  $R(\cdot)$  is a deterministic realization of  $\mathcal{R}(\cdot)$  and  $w \in \mathbb{R}^d, x: \mathcal{S} \rightarrow \mathbb{R}^d$  represent the weight and the feature. The IRL problem is to identify the weight vector  $w$ , given a set of demonstrations. The resulting value function for a policy  $\pi$  can be expressed as

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right] = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t w^T x(s_t) \mid s_0 = s \right] \\ &= w^T \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t x(s_t) \mid s_0 = s \right] = w^T \mu(\pi), \end{aligned}$$

where  $\mu(\pi \mid s_0 = s) \in \mathbb{R}^d$  is the discounted weighted frequency of state features  $x(s)$  under policy  $\pi$ .



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen

# Inverse RL via Margin Optimization

$$\mathbb{E}_{\pi^*} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid s_0 = s \right] \geq \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid s_0 = s \right], \quad \forall \pi,$$

where  $R^*$  denotes an optimal reward function. Thus, if an expert's demonstrations are optimal (i.e. actions are drawn from an optimal policy), **to identify  $w$  it is sufficient to find some  $w^*$**  such that

$$w^{*T} \mu(\pi^* \mid s_0 = s) \geq w^{*T} \mu(\pi \mid s_0 = s), \quad \forall \pi, \forall s,$$

where some restrictions are put on  $w^*$  to avoid trivial solutions to the linear system. As long as this constraint is linear, the problem can be solved by linear programming.

$$\max_{w^{*T}} w^{*T} \mu(\pi^* \mid s_0 = s) - w^{*T} \mu(\pi \mid s_0 = s), \quad \forall \pi^* \neq \pi, \forall s,$$

$$s.t., \|w^{*T}\| = 1$$



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen

# Inverse RL via Maximum entropy

There is an infinite number of reward functions with the same optimal policy, and an infinite number of stochastic policies that can match feature counts. To address the problem of ambiguity, **Maximum Entropy (MaxEnt) IRL** considers the collection of all possible  $H$ -step trajectories in a deterministic MDP. For a linear reward model, **a policy is completely specified by its distribution over trajectories.**

*Given this, which policy should we choose given a set of  $k$  distributions?*



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Inverse RL via Maximum entropy

The principle of maximum entropy motivates choosing a distribution with **no additional** preferences beyond matching the feature expectations in the demonstration dataset

$$\text{Maximize } - \sum_{\tau} p(\tau) \log p(\tau)$$

$$\text{Subject to } \int p(\tau) r(\tau) d\tau = \frac{1}{N} \sum_{\tau \in D_E} r(\tau), \text{ and } \int p(\tau) d\tau = 1$$

where  $D_E$  denotes the expert dataset. This is equivalent to specifying the **reward function**  $r$  that yields a policy with the maximum entropy, constrained to matching the **reward expectations** in the expert demonstration.



# Inverse RL via Maximum entropy

Maximizing the entropy of the distribution over the paths subject to the reward constraints from observed data implies we maximize the likelihood of the observed data under the maximum entropy (exponential family) distribution

$$p(D_E|r_w) = \prod_{i=1}^N p(\tau^{(i)}|r_w) = \frac{1}{(Z_M)^N} \prod_{i=1}^N \exp[r_w(\tau^{(i)})] \quad (1)$$

where the normalizing term  $Z_M = \int \exp[r_w(\tau)] d\tau$ .





# Inverse RL via Maximum entropy

The gradient is the difference between expected empirical rewards and the learner's expected rewards, which can be expressed as follows

$$\nabla_w \log [p(D_E | r_w)] = \sum_{i=1}^N \left[ \nabla_w \sum_{t=0}^T \log[r_w(s_t^{(i)}, a_t^{(i)})] \right] - N E_{\hat{\tau} \sim p(\tau | r_w)} \left[ \nabla_w \sum_{t=0}^T \log[r_w(\hat{s}_t, \hat{a}_t)] \right]$$



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Question and Answering (Q&A)



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen