

CSE 571
Imitation Learning and Policy gradient
 Dieter Fox

Slides borrowed from many sources -
 Abhishek Gupta, Liyiming Ke, Sergey Levine

1

How can we learn policies?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

Model-free RL

Model-based RL

Imitation Learning

2

Idea 1: Imitation Learning via Behavior Cloning

Given: Demonstrations of optimal behavior $\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)]$
 Goal: Train a policy to mimic the demonstrator

Idea: Treat imitation learning as a supervised learning problem!

3

So does behavior cloning really work?

- Imitation Learning \neq Supervised Learning

$$\arg \max_{\theta} \mathbb{E}_{(s^*, a^*) \sim \mathcal{D}} [\log \pi_{\theta}(a^* | s^*)] \quad \mathbb{E}_{(s, a) \sim \rho(\pi)} [1(a = a^*)]$$

Not the same!

4

What is the general principle?

— training trajectory
— π_θ expected trajectory

stability

Corrective labels that bring you back to the data

5

Concrete Instantiation: DAgger

can we make $p_{\text{data}}(\mathbf{o}_t) = p_{\pi_\theta}(\mathbf{o}_t)$?
 idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$
 how? just run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$
 but need labels \mathbf{a}_t !

1. train $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
2. run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

Ross et al. '11

6

Why might we fail to fit the expert?

Multimodal behavior.. amongst other reasons

Not a matter of network size! It's about distributional expressivity

7

Why might we fail to fit the expert?

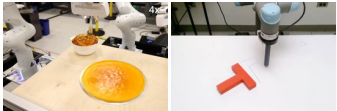
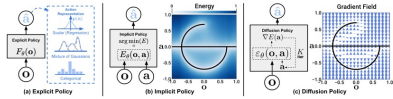
Multimodal behavior → use more expressive probability distributions

1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...

8


Why might we fail to fit the expert?

1. Output mixture of Gaussians
2. Latent variable models
3. Autoregressive discretization
4. Diffusion models
5. ...

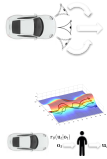



9

Perspectives on Imitation – don't believe everything you see online



- Pros:
 - Easy to use, no additional infra
 - Can sometimes be unreasonably effective
- Cons:
 - Challenges of compounding error, multimodality
 - Doesn't really generalize
 - Very expensive in terms of data collection!

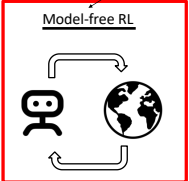


10


How can we learn policies?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$


Model-free RL



Model-based RL



Imitation Learning



11

What if we just performed gradient ascent?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right] = \int p_{\theta}(\tau) R(\tau) d\tau$$

Standard gradient descent (supervised learning)

$$\nabla_{\theta} \mathbb{E}_{x \sim g(x)} [f_{\theta}(x)]$$

REINFORCE gradient descent (RL)

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}(x)} [f(x)]$$

12

Taking the gradient of return

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\nabla_{\theta} \log p_{\theta}(\tau) \sum_{t=0}^T r(s_t, a_t) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t) \\ a_t \sim \pi(a_t|s_t)}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \sum_{t'=0}^T r(s_{t'}, a_{t'}) \right]$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i|s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i) \quad (\text{approximating using samples})$$

13

What does this mean?

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i|s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Increase the likelihood of actions in high return trajectory:

14

Resulting Algorithm (REINFORCE)

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

Collect Data

Take Gradient Step

On-policy → REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(a_t|s_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i|s_t^i)) (\sum_{t'} r(s_{t'}^i, a_{t'}^i))$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

15

Policy Gradient in Action

Iteration 0

16