



Introduction to Imitation Learning

Anup Deshmukh

Advisor: Prof. Dinesh Babu

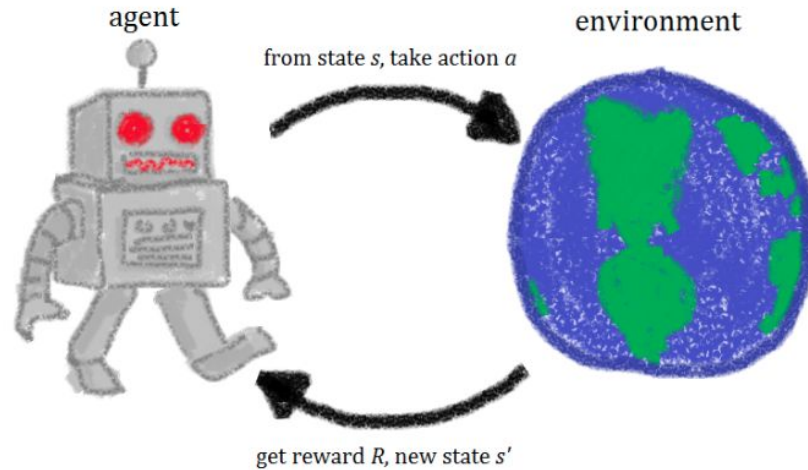
IIIT-Bangalore



Sections

- Reinforcement Learning 101
- Why are we not THERE yet?
- What is Imitation Learning?
- **Challenges in Imitation Learning!**
- **Posed questions and Ideas.**

1 Reinforcement Learning 101



Usual Reinforcement Learning setup



1 Reinforcement Learning 101: Approach 1

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s]$$

Expected Reward discounted Given that state

Value based : The value function tells us the maximum expected future reward the agent will get at each state.

Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

TD Learning

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Previous estimate Reward t+1 Discounted value on the next step TD Target



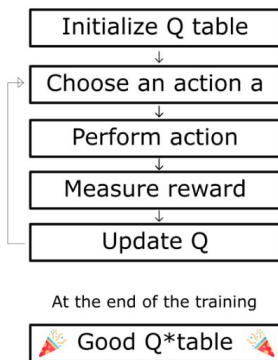
1 Reinforcement Learning 101: Approach 2

$$\pi_{\theta}(a|s) = P[a|s]$$

Probability of taking action a given state s with parameters theta.

Policy based: In policy-based RL, we want to directly optimize the policy function $\pi(s)$ without using a value function.

1 Reinforcement Learning 101: Learning Method 1



Q learning: We learn the action value function by maintaining the Q table



	←	→	↑	↓
Start	0	0	0	0
Small cheese	0	0	0	0
Nothing	0	0	0	0
2 small cheese	0	0	0	0
Death	0	0	0	0
Big cheese	0	0	0	0

The initialized Q-table

1. Initialize Q-values ($Q(s, a)$) arbitrarily for all state-action pairs.
2. For life or until learning is stopped...
3. Choose an action (a) in the current world state (s) based on current Q-value estimates ($Q(s, \cdot)$).
4. Take the action (a) and observe the outcome state (s') and reward (r).
5. Update $Q(s, a) := Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$

1 Reinforcement Learning 101: Learning Method 2

$$\pi_{\theta}(a|s) = P[a|s]$$

Probability of taking action a given state s with parameters theta.

$$J_{avgv}(\theta) = E_{\pi}(V(s)) = \sum d(s)V(s)$$

$$\text{where } d(s) = \frac{N(s)}{\sum_{s'} N(s')}$$

Number of occurrences of the state

Total nb occurrences of all states

$$J_{avR}(\theta) = E_{\pi}(r) = \sum_s d(s) \sum_a \pi_{\theta}(s, a) R_s^a$$

Probability that I'm in state s

Probability that I take this action a from that state under this policy

Immediate reward that I'll get

Policy gradient: In policy-based RL, we want to directly optimize the policy function $\pi(s)$ without using a value function.

Policy : π_{θ}

Objective function : $J(\theta)$

Gradient : $\nabla_{\theta} J(\theta)$

Update : $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

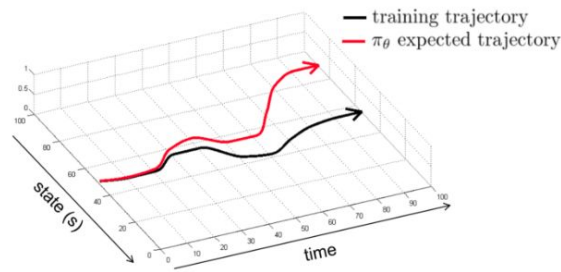


2 Problems in RL

- Sparsity in Rewards
 - How do you know when you have made a catastrophic move?
 - Take an example of game of Chess or Recommender Systems
- Sample inefficiency
 - How do you speed up learning in real life?
 - Again take an example of Recommender Systems

3 What is Imitation Learning?

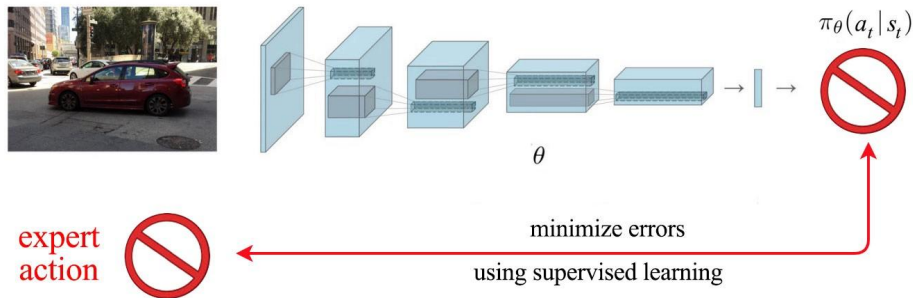
- Idea is to implicitly give an agent prior information about the world by mimicking human behavior.
- Pretraining it on a human demonstrator's data might also make the training process faster



Source

Imitation Learning





Algorithm 43 `SUPERVISEDIMITATIONTRAIN`($\mathcal{A}, \tau_1, \tau_2, \dots, \tau_N$)

1: $D \leftarrow \langle (x, a) : \forall n, \forall (x, a, \ell) \in \tau_n \rangle$ // collect all observation/action pairs
 2: **return** $\mathcal{A}(D)$ // train multiclass classifier on D

Algorithm 44 `SUPERVISEDIMITATIONTEST`(f)

1: **for** $t = 1 \dots T$ **do**
 2: $x_t \leftarrow$ current observation
 3: $a_t \leftarrow f(x_t)$ // ask policy to choose an action
 4: take action a_t
 5: $\ell_t \leftarrow$ observe instantaneous loss
 6: **end for**
 7: **return** $\sum_{t=1}^T \ell_t$ // return total loss



Imitation vs. Reinforcement Learning

imitation learning

- Requires demonstrations
- Must address distributional shift
- Simple, stable supervised learning
- Only as good as the demo

reinforcement learning

- Requires reward function
- Must address exploration
- Potentially non-convergent RL
- Can become arbitrarily good

Can we get the best of both?

e.g., what if we have demonstrations *and* rewards?



Simplest combination: pretrain & finetune

- Demonstrations can overcome exploration: show us how to do the task
- Reinforcement learning can improve beyond performance of the demonstrator
- Idea: initialize with imitation learning, then finetune with reinforcement learning!

1. collected demonstration data $(\mathbf{s}_i, \mathbf{a}_i)$
2. initialize π_θ as $\max_\theta \sum_i \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i)$




3. run π_θ to collect experience
4. improve π_θ with any RL algorithm




Simplest combination: pretrain & finetune

Pretrain & finetune

1. collected demonstration data $(\mathbf{s}_i, \mathbf{a}_i)$
2. initialize π_θ as $\max_\theta \sum_i \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i)$
-  3. run π_θ to collect experience
4. improve π_θ with any RL algorithm

vs. DAgger

What will you do when you drift
off-course?

- 
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: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

5 Connection to RS: Why do we need policy gradient?

$$\pi_{\theta}(a|s) = P[a|s]$$

Probability of taking action a given state s with parameters theta.

$$J_{avR}(\theta) = E_{\pi}(r) = \sum_s d(s) \sum_a \pi_{\theta}(s, a) R_s^a$$

Probability that I'm in state s

Probability that I take this action a from that state under this policy

Immediate reward that I'll get

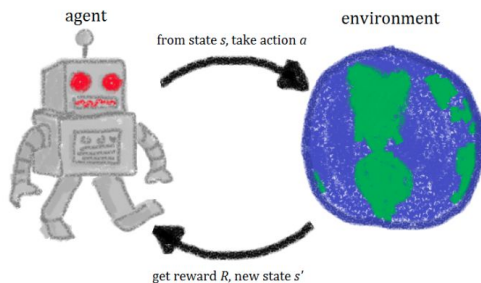
Policy : π_{θ}
Objective function : $J(\theta)$
Gradient : $\nabla_{\theta} J(\theta)$
Update : $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Can we arrive at the optimal policy for the generator in GAN's using the Imitation Learning?

AND

Do away with this optimization?

5 Connection to RS



Usual Reinforcement Learning setup

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t \underline{r(s_t, a_t)} \right]$$

$$\underbrace{p_{\theta}(s_1, a_1, \dots, s_T, a_T)}_{\tau} = p(s_1) \prod_{t=1}^T \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}} \underbrace{p(s_{t+1} | s_t, a_t)}_{\text{model}}$$

$$\underbrace{p_{\theta}(s_1, a_1, \dots, s_T, a_T)}_{\text{demonstration } p(\tau)} = \underbrace{p(s_1)}_{\text{train}} \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$



Thank
you.

