

Last Time

- What tools do we have to solve MDPs with continuous S and A ?

Course Map

- Outcome Uncertainty, Immediate vs Future Rewards (MDP)
- Model Uncertainty (Reinforcement Learning)
- State Uncertainty (POMDP)
- Interaction Uncertainty (Game)

Course Map

- Outcome Uncertainty, Immediate vs Future Rewards (MDP)
- Model Uncertainty (Reinforcement Learning)
- State Uncertainty (POMDP)
- Interaction Uncertainty (Game)



Course Map

- Outcome Uncertainty, Immediate vs Future Rewards (MDP)
- Model Uncertainty (Reinforcement Learning)
- State Uncertainty (POMDP)
- Interaction Uncertainty (Game)



Guiding Questions

Guiding Questions

- What is Reinforcement Learning?
- What are the main challenges in Reinforcement Learning?

Guiding Questions

- What is Reinforcement Learning?
- What are the main challenges in Reinforcement Learning?
- How do we categorize RL approaches?

Problem from HW2

Question 2. (25 pts) Consider a game with 3 squares in a horizontal line drawn on paper, a token, and a die. Each turn, the player can either reset or roll the die. If the player rolls and the die shows an odd number, the token is moved one square to the right, and if an even number is rolled, the token is moved two squares to the right (in both cases stopping at the rightmost square¹). If the player resets, the token is always moved to the leftmost square. If the reset occurs when the token is in the middle square, two points are added; if the player resets when the token is on the right square, a point is subtracted.

- c) Suppose you are not sure that the die is fair (i.e. whether it will yield odd and even with equal probability). Give finite upper and lower bounds for the accumulated discounted score that you can expect to receive with discount $\gamma = 0.95$.

$$\frac{1-\gamma}{1-\gamma}$$

Reinforcement Learning

Reinforcement Learning

Previously: (S, A, T, R, γ)

Reinforcement Learning

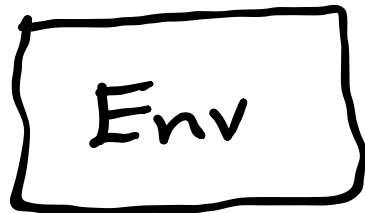
Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$

Unknown!

Reinforcement Learning

Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$
Unknown!

Now: Episodic Simulator

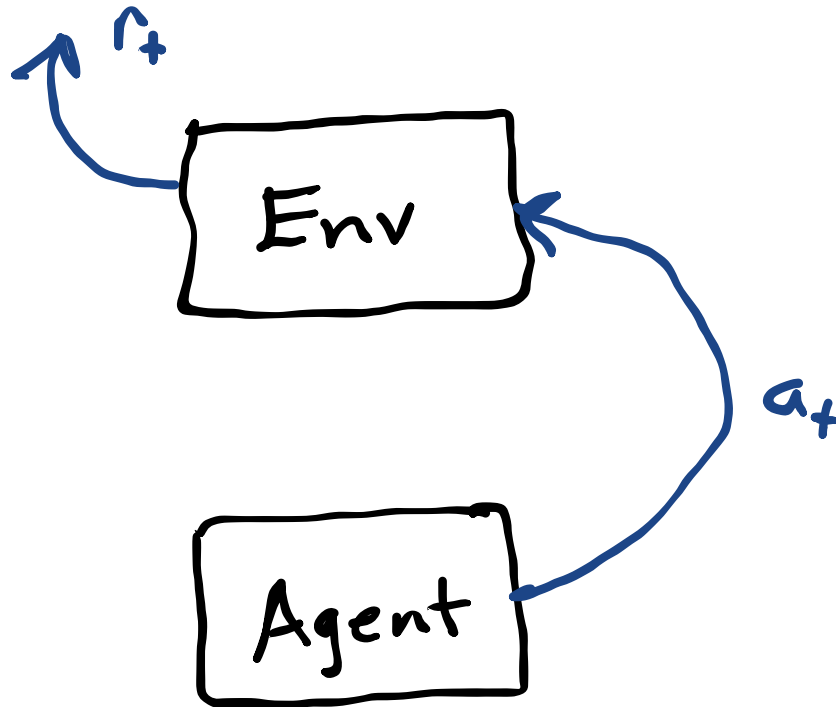
A hand-drawn rectangular box with a slightly irregular border, containing the text "Env" in a handwritten style.A hand-drawn rectangular box with a slightly irregular border, containing the text "Agent" in a handwritten style.

Reinforcement Learning

Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$

Unknown!

Now: Episodic Simulator



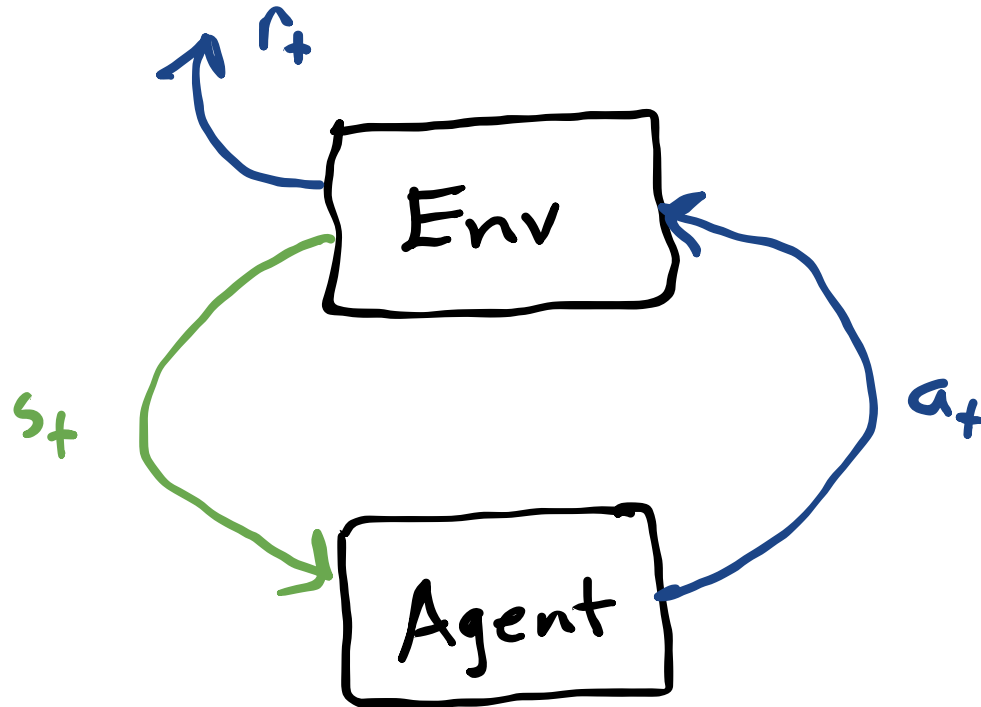
$r = \text{act!}(\text{env}, a)$

Reinforcement Learning

Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$

Unknown!

Now: Episodic Simulator



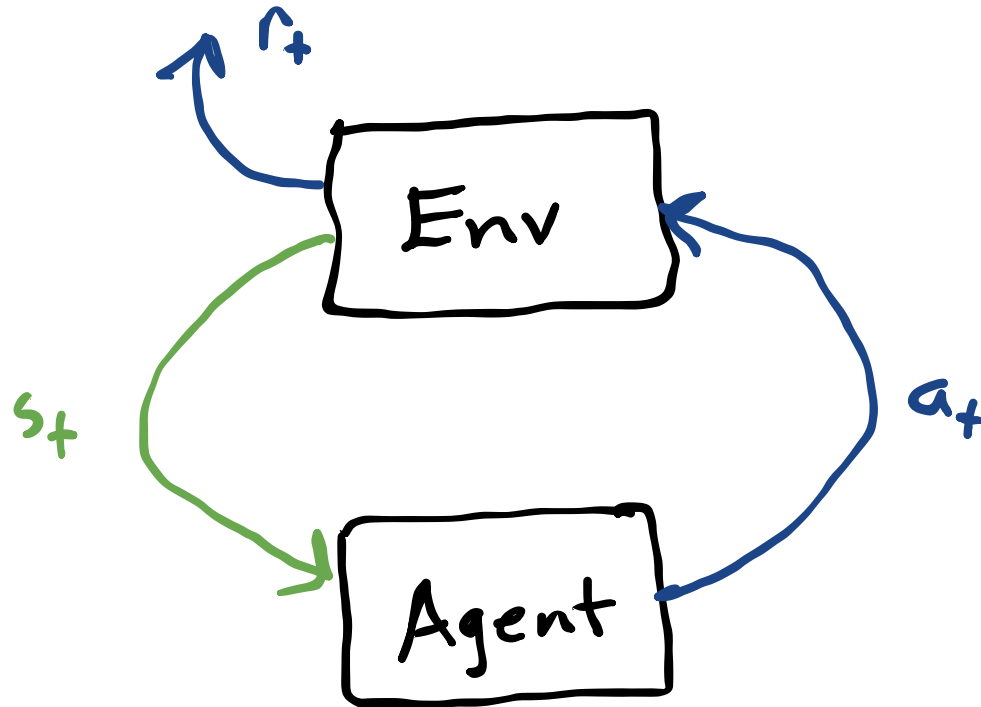
$r = \text{act!}(\text{env}, a)$

$s = \text{observe}(\text{env})$

Reinforcement Learning

Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$
Unknown!

Now: Episodic Simulator



`r = act!(env, a)`

`s = observe(env)`

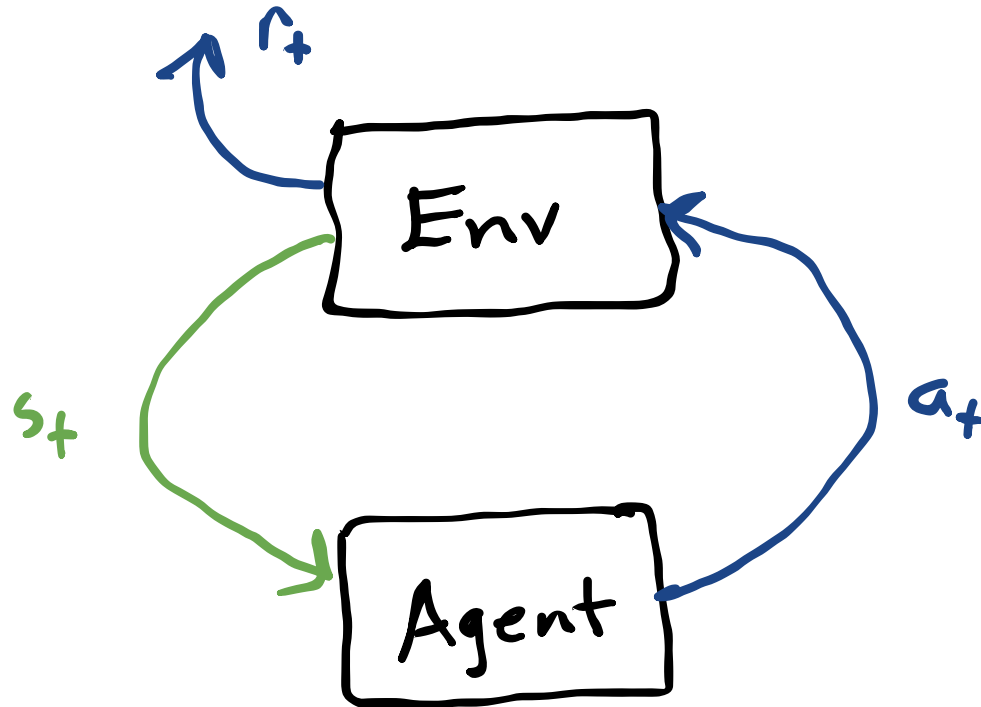
In python, typically
`s, r = env.step(a)`

Reinforcement Learning

Previously: $(S, A, \cancel{T}, \cancel{R}, \gamma)$

Unknown!

Now: Episodic Simulator



`r = act!(env, a)`

`s = observe(env)`

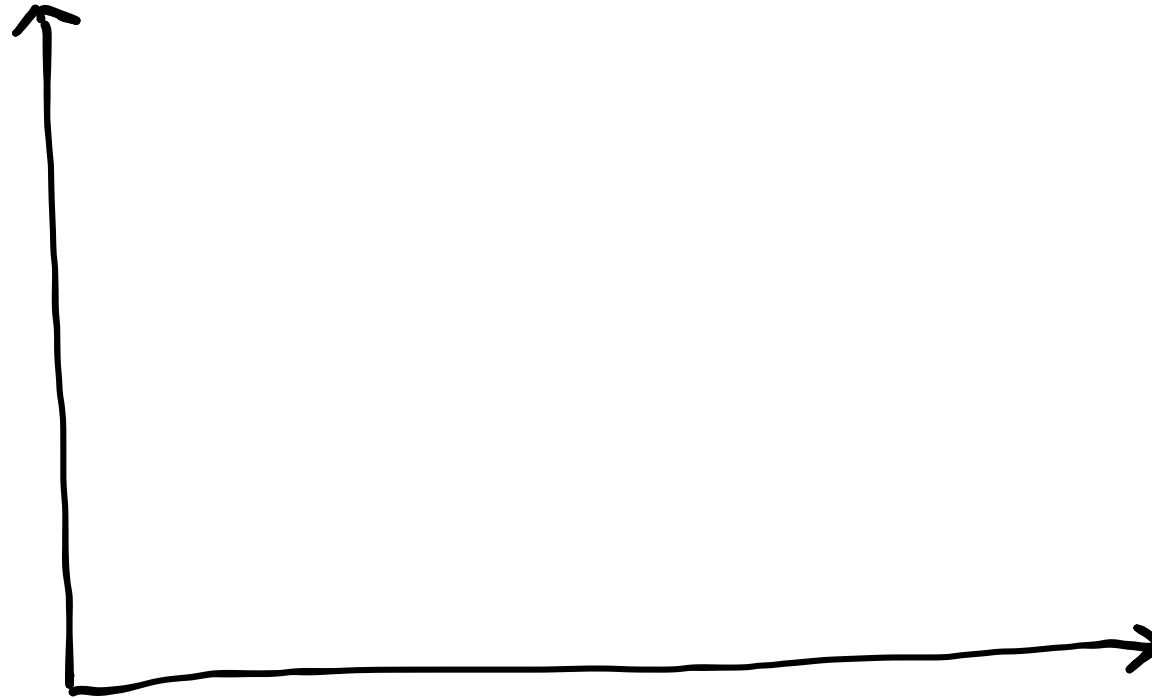
In python, typically

`s, r = env.step(a)`

Note: Different from $s', r = G(s, a)$

Learning Curve

Learning Curve



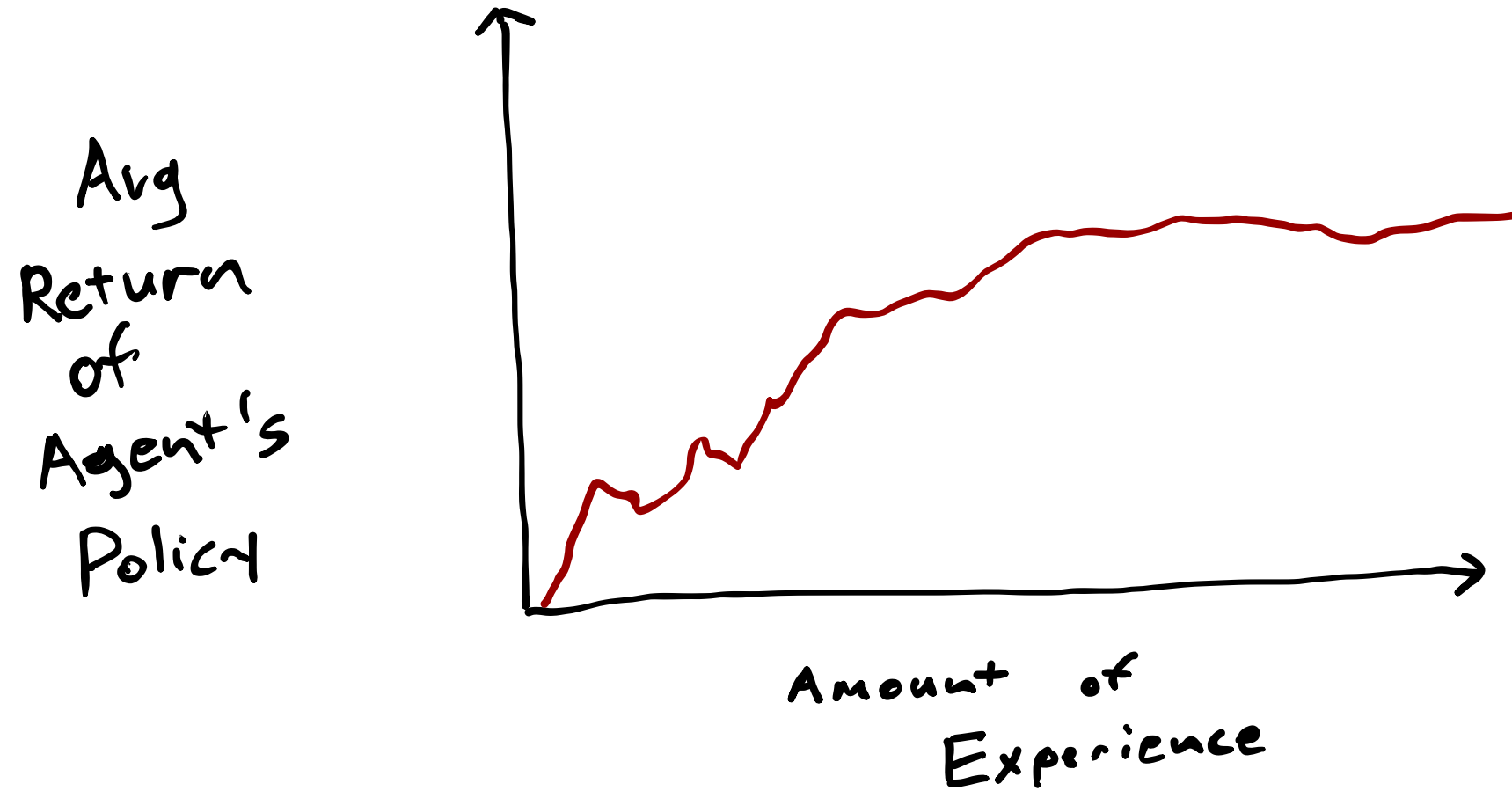
Learning Curve



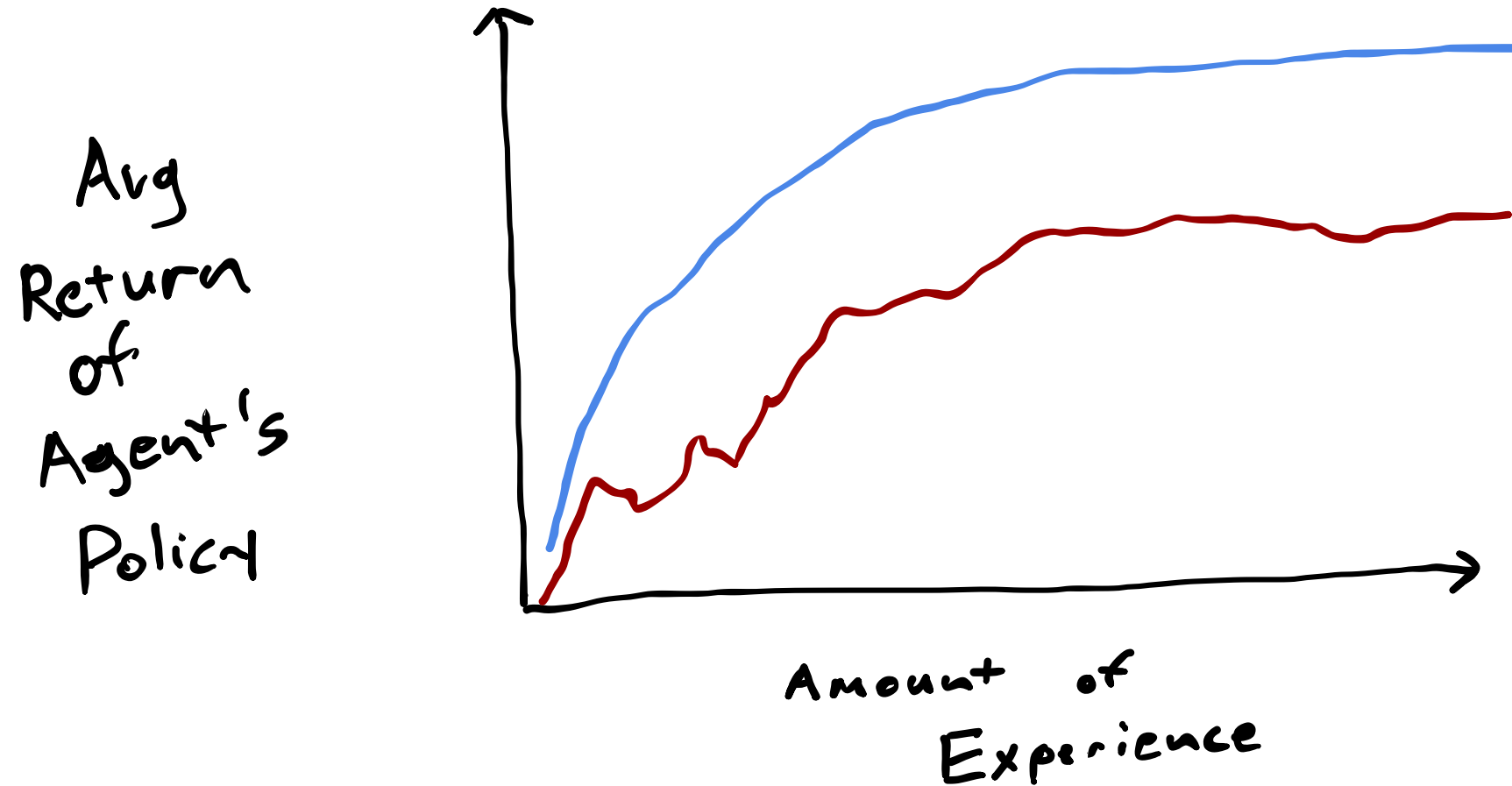
Learning Curve



Learning Curve



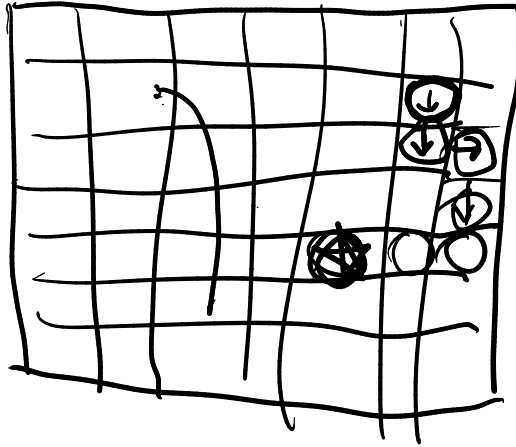
Learning Curve



Break

know S, A

How to interact to maximize reward
given no knowledge of T and R .



$$\textcircled{1} \textcircled{\hat{R}(s,a)} = \left\{ \begin{array}{l} \hat{U} \\ \hat{Q}(s,a) \\ \hat{\pi}(s) \end{array} \right.$$

heuristic that steers toward
reward states,

$$\pi(s) = \underset{a}{\operatorname{argmax}} \left(\hat{R}(s,a) + \gamma \textcircled{E[\hat{U}(s')]} \right)$$

$$\pi(s) = \underset{a}{\operatorname{argmax}} \hat{Q}(s,a)$$

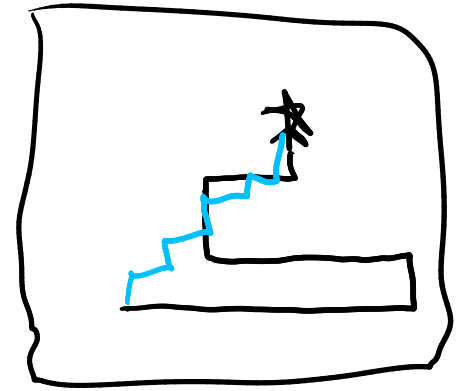
$$\hat{\pi}(s)$$

= action that has led to the most future
reward in previous
episodes₇

Challenges

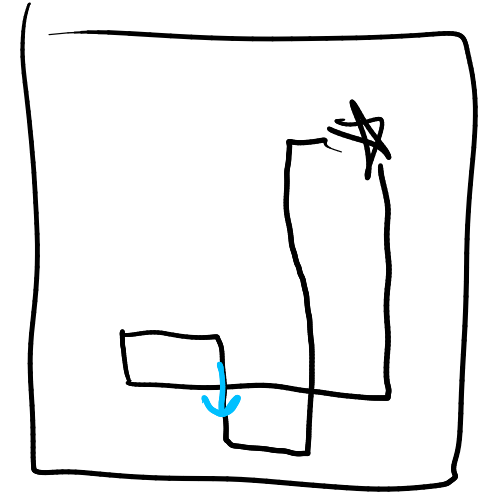
Challenges

1. Exploration vs Exploitation



Challenges

1. Exploration vs Exploitation
2. Credit Assignment



Challenges

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Classifications

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
- **Model Free:** Attempt to find Q^* or π^* directly without estimating T or R

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
- **Model Free:** Attempt to find Q^* or π^* directly without estimating T or R
- **On-Policy:** Learn only using experience generated with the current policy.

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
- **Model Free:** Attempt to find Q^* or π^* directly without estimating T or R
- **On-Policy:** Learn only using experience generated with the current policy.
- **Off-Policy:** Learn using experience generated from the current policy *and* previous policies.

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
- **Model Free:** Attempt to find Q^* or π^* directly without estimating T or R
- **On-Policy:** Learn only using experience generated with the current policy.
- **Off-Policy:** Learn using experience generated from the current policy *and* previous policies.
- **Batch:** Learn only from previously-generated experience.

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
- **Model Free:** Attempt to find Q^* or π^* directly without estimating T or R
- **On-Policy:** Learn only using experience generated with the current policy.
- **Off-Policy:** Learn using experience generated from the current policy *and* previous policies.
- **Batch:** Learn only from previously-generated experience.
 - **Tabular:** Keep track of learned values for each state in a table

Classifications

- • **Model Based**: Attempt to learn T and R , then find π^* by solving MDP
- **Model Free**: Attempt to find Q^* or π^* directly without estimating T or R
- **On-Policy**: Learn only using experience generated with the current policy.
- • **Off-Policy**: Learn using experience generated from the current policy *and* previous policies.
- **Batch**: Learn only from previously-generated experience.
- • **Tabular**: Keep track of learned values for each state in a table
- **Deep**: Use a neural network to approximate learned values

Tabular Maximum Likelihood Model-Based RL

TMLMBRL
MLMBTRL

Given env, S, A

$N \leftarrow 0$

$\rho \leftarrow 0$

$s \leftarrow \text{observe}(env)$

$\pi \leftarrow \text{random policy}$

loop

$a \leftarrow \begin{cases} \text{rand}(A) & \text{w.p. } \epsilon \\ \pi(s) & \text{w.p. } 1-\epsilon \end{cases}$

$r \leftarrow \text{act!}(env, a)$

$s' \leftarrow \text{observe}(env)$

$N[s, a, s'] += 1$

$\rho[s, a] += r$

$T[s, a, s'] \leftarrow \frac{N[s, a, s']}{\sum_{s'} N[s, a, s']}$

$R[s, a] \leftarrow \frac{\rho[s, a]}{\sum_{s'} N[s, a, s']}$

$\pi \leftarrow \text{solve}(T, R)$

$s \leftarrow s'$

$N[s, a, s'] \leftarrow \# \text{ a taken in } s \text{ resulting in } s'$
 $\rho[s, a] \leftarrow \text{cumulative reward}$

(s, a, r, s')

$\forall s, a, s'$

$\forall s, a$

← expensive

Guiding Questions

- What is Reinforcement Learning?
- What are the main challenges in Reinforcement Learning?
- How do we categorize RL approaches?

*Exploitation vs Exploration
Credit Assignment
Gen.*