

Reinforcement Learning

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 $\gamma = 0.95$.

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

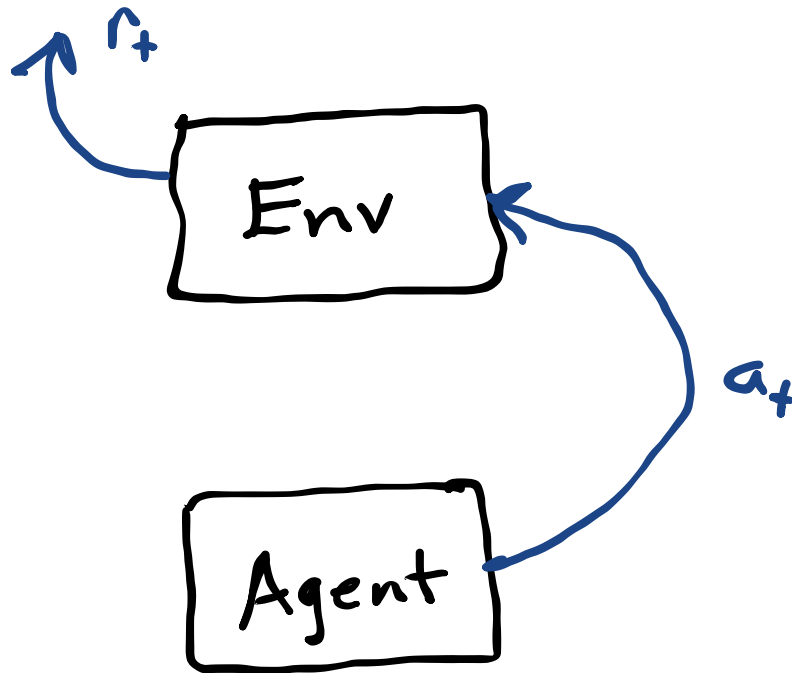
Env

Agent

Reinforcement Learning

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

Now: Episodic Simulator

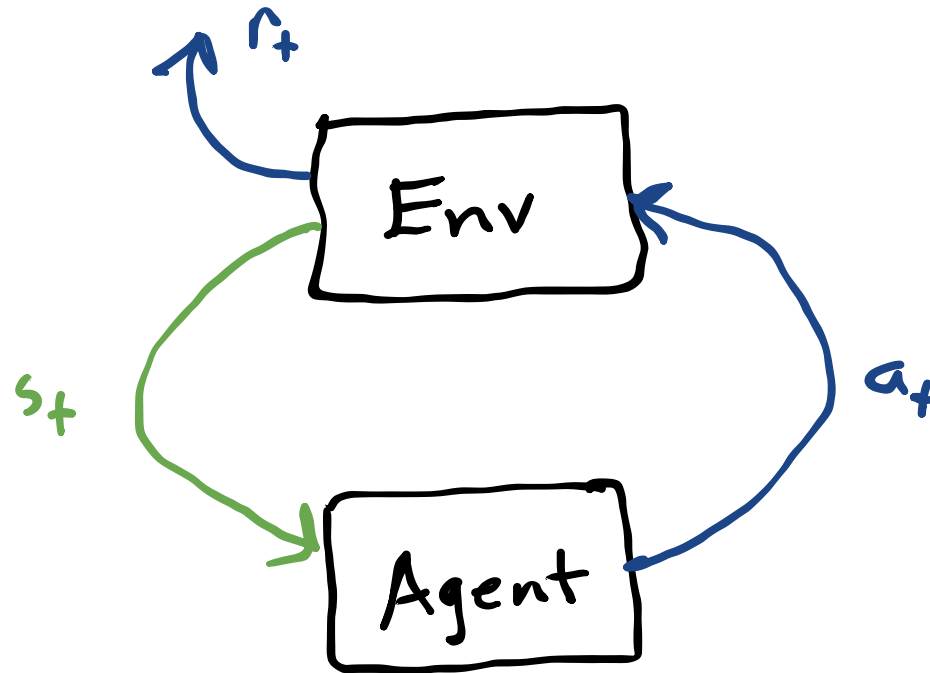


$r = \text{act!}(\text{env}, a)$

Reinforcement Learning

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

Now: Episodic Simulator



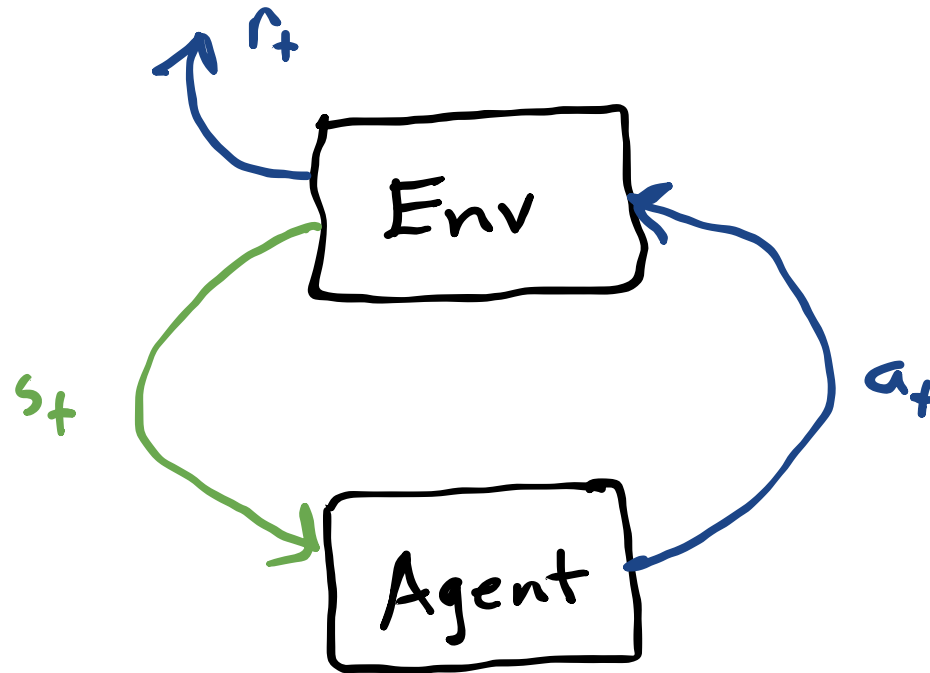
$r = \text{act!}(\text{env}, a)$

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

Reinforcement Learning

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

Now: Episodic Simulator



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

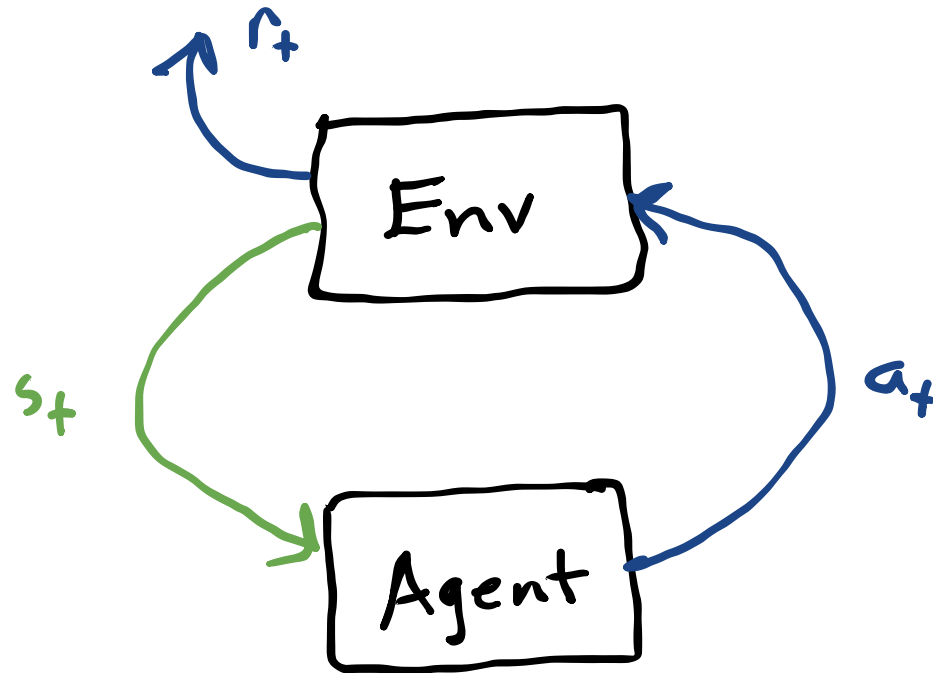
$s = \text{observe}(\text{env})$

In python, typically
 $s, r = \text{env.step}(a)$

Reinforcement Learning

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

Now: Episodic Simulator



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

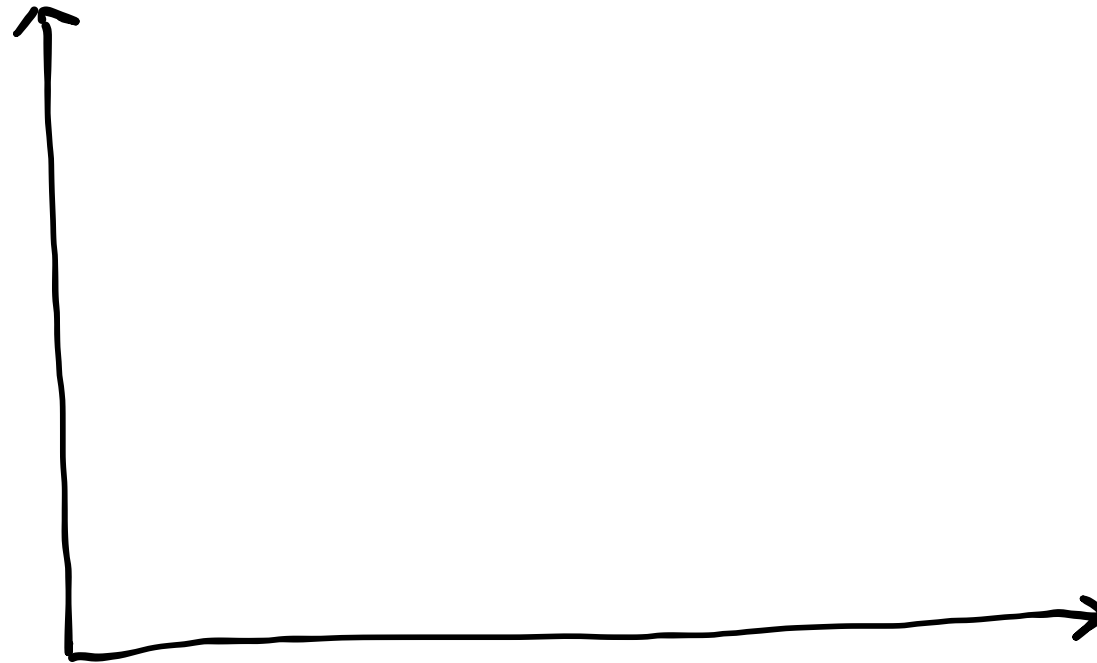
$s = \text{observe}(\text{env})$

In python, typically
 $s, r = \text{env.step}(a)$

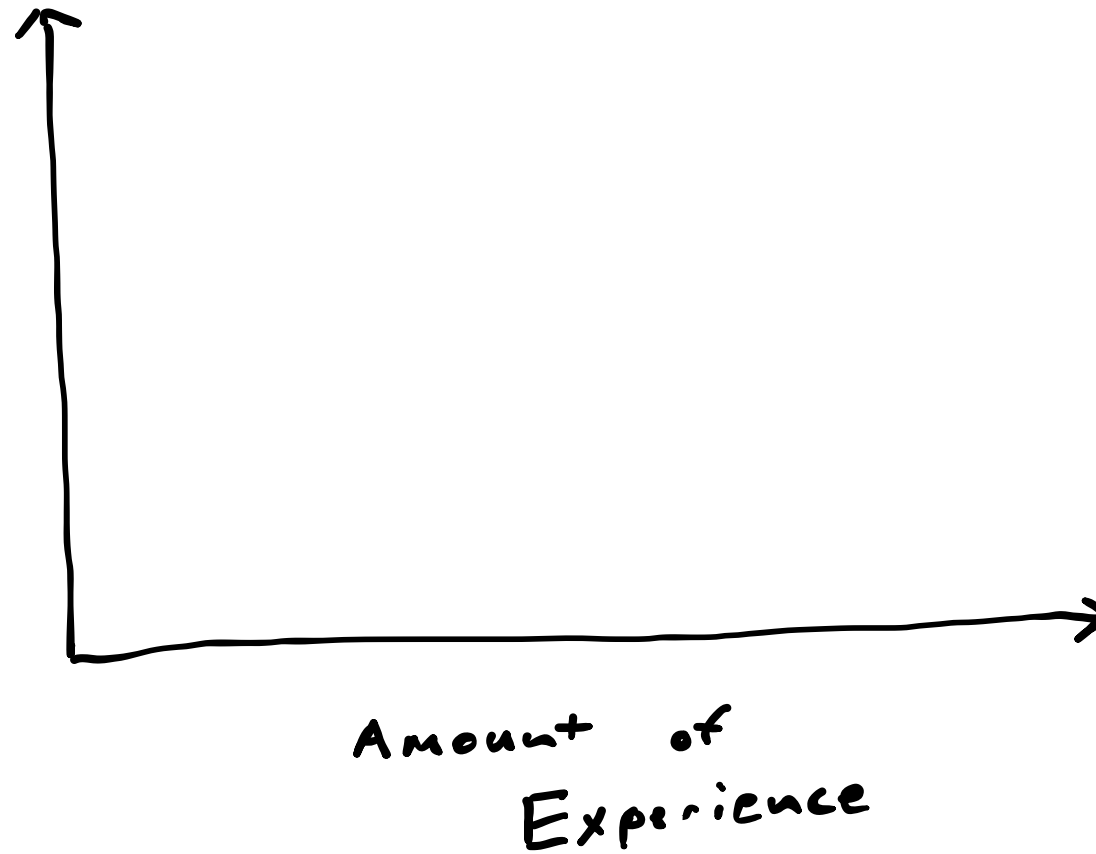
Note: Different from $s', r = G(s, a)$

Learning Curve

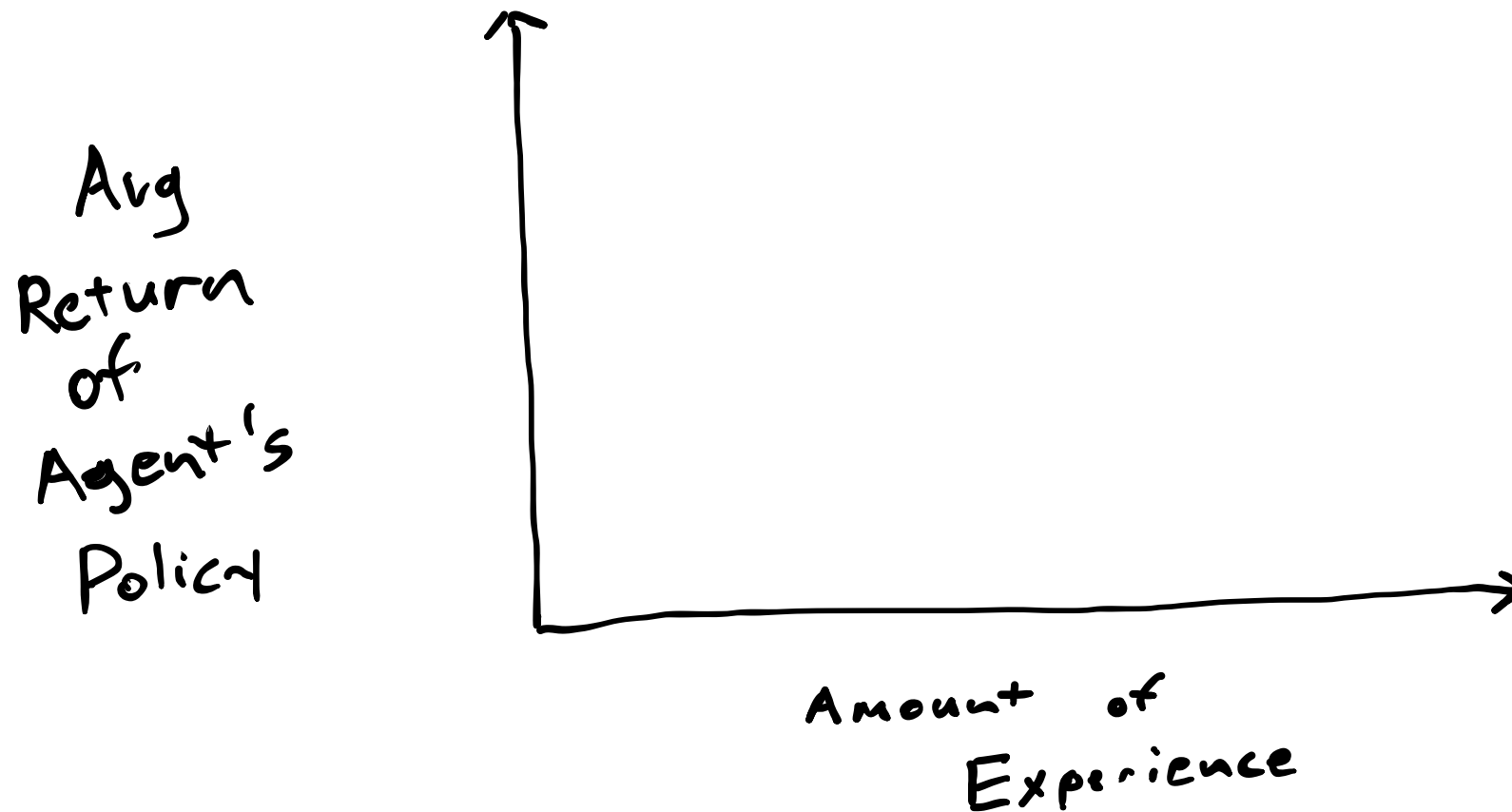
Learning Curve



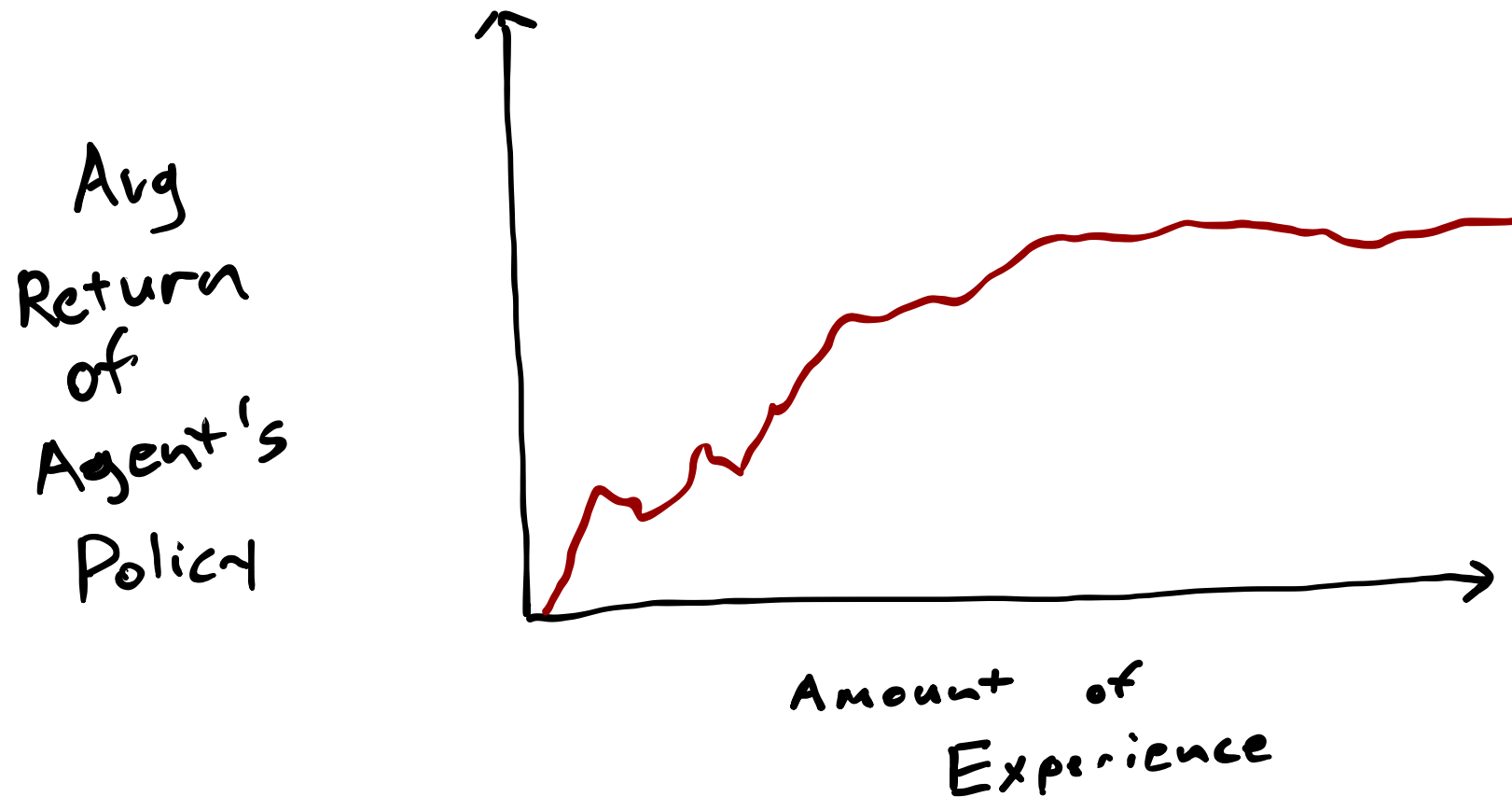
Learning Curve



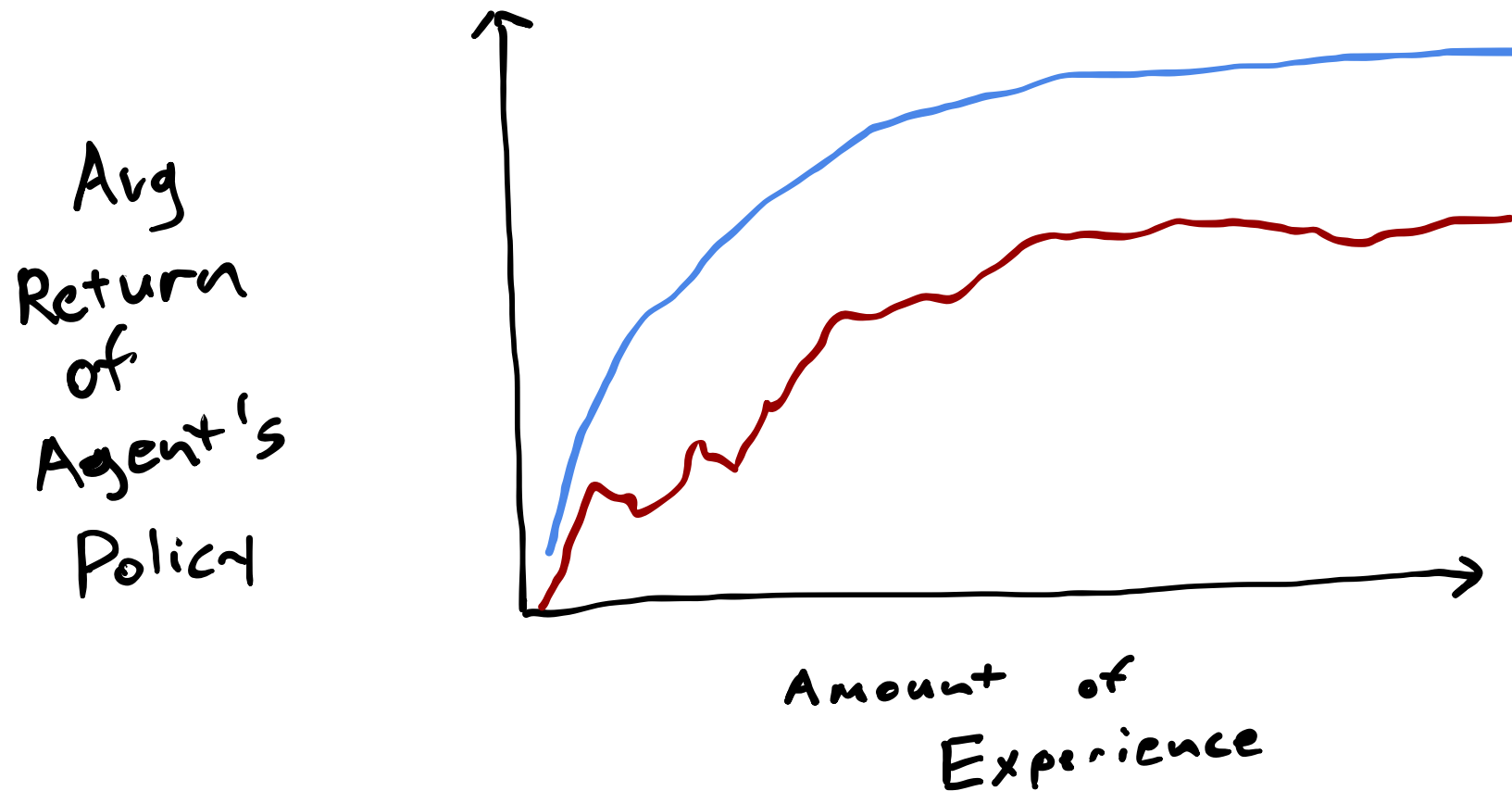
Learning Curve



Learning Curve



Learning Curve



Break

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

Classifications

- **Model Based:** Attempt to learn T and R , then find π^* by solving MDP
 - **Model Free:** Attempt to find Q^* or π^* directly
-
- **On-Policy:** The exploration policy is the same as the learned 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
-
- **On-Policy:** The exploration policy is the same as the learned policy.
 - **Off-Policy:** The exploration policy may be different than the learned 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
-
- **On-Policy:** The exploration policy is the same as the learned policy.
 - **Off-Policy:** The exploration policy may be different than the learned policy.
 - **Batch:** Learn only from previously-generated experience (no exploration 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
-
- **On-Policy:** The exploration policy is the same as the learned policy.
 - **Off-Policy:** The exploration policy may be different than the learned policy.
 - **Batch:** Learn only from previously-generated experience (no exploration policy).
-
- **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
-
- **On-Policy:** The exploration policy is the same as the learned policy.
 - **Off-Policy:** The exploration policy may be different than the learned policy.
 - **Batch:** Learn only from previously-generated experience (no exploration policy).
-
- **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

Given env, S , A

$N[s, a, s'] \leftarrow 0 \quad \forall s, a, s'$

$\rho[s, a] \leftarrow 0 \quad \forall s, a$

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

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

loop

 reset!(env)

while not terminated(env)

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

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

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

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

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

$s \leftarrow s'$

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

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

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

Guiding Questions

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