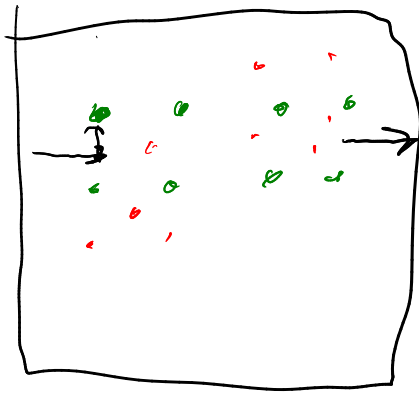


1

# Reinforcement Learning



heuristic-pol( $m, s$ )

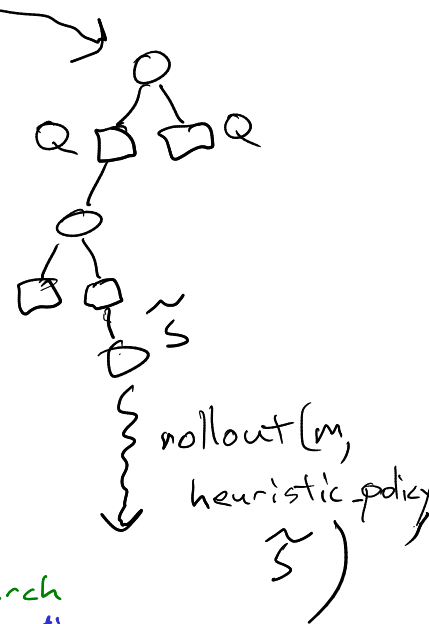
for a  
 $s', r \leftarrow @gen(sp, r)(m, sa)$

rollout( $m, select\_action,$  )

for t

$a = select\_action(m, s)$

$sp, r = @gen(sp, r)(m, s, a)$

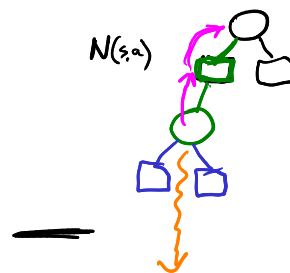


Algorithm 4.9 Monte Carlo tree search

```

1: function SELECTACTION( $s, d$ )
2:   loop
3:     SIMULATE( $s, d, \pi_0$ )
4:   return  $\arg \max_a Q(s, a)$ 
5: function SIMULATE( $s, d, \pi_0$ )
6:   if  $d = 0$ 
7:     return 0
8:   if  $s \notin T$ 
9:     for  $a \in A(s)$ 
10:      ( $N(s, a), Q(s, a)$ )  $\leftarrow$  ( $N_0(s, a), Q_0(s, a)$ )
11:     $T = T \cup \{s\}$ 
12:    return ROLLOUT( $s, d, \pi_0$ )
13:    $a \leftarrow \arg \max_{a \in A(s)} \left[ Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}} \right]$ 
14:   ( $s', r$ )  $\sim G(s, a)$ 
15:    $q \leftarrow r + \gamma \text{SIMULATE}(s', d-1, \pi_0)$ 
16:    $N(s, a) \leftarrow N(s, a) + 1$ 
17:    $Q(s, a) \leftarrow Q(s, a) + \frac{q - Q(s, a)}{N(s, a)}$ 
18:   return  $q$ 

```



search  
 expansion  
 Value Estimate/rollout  
 Backup

simulate! ( $m, s, n, q, t = \text{nothing}$ )

$\gamma^{100}$

Algorithm 4.10 Rollout evaluation

```

1: function ROLLOUT( $s, d, \pi_0$ )
2:   if  $d = 0$ 
3:     return 0
4:    $a \sim \pi_0(s)$ 
5:   ( $s', r$ )  $\sim G(s, a)$ 
6:   return  $r + \gamma \text{ROLLOUT}(s', d-1, \pi_0)$ 

```

# Last Time

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

- LQR  $\pi^*(s) = -Ks$
- Sparse Sampling / PW in tree search
- Value Func. Approx.
- MPC



# 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 square<sup>1</sup>). 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:  $(S, A, T, R, \gamma)$

# 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

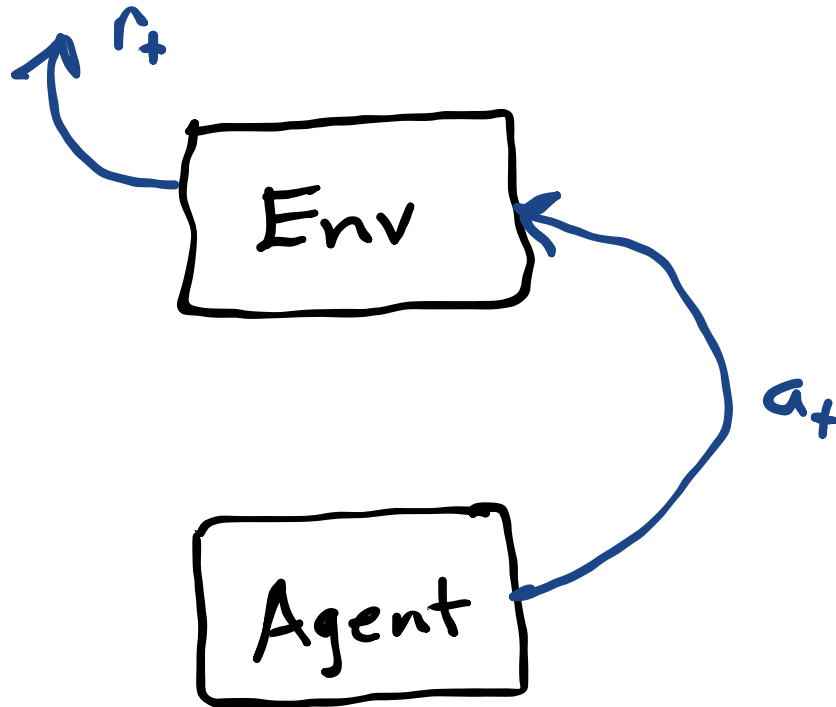
Agent

# Reinforcement Learning

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

Unknown!

Now: Episodic Simulator



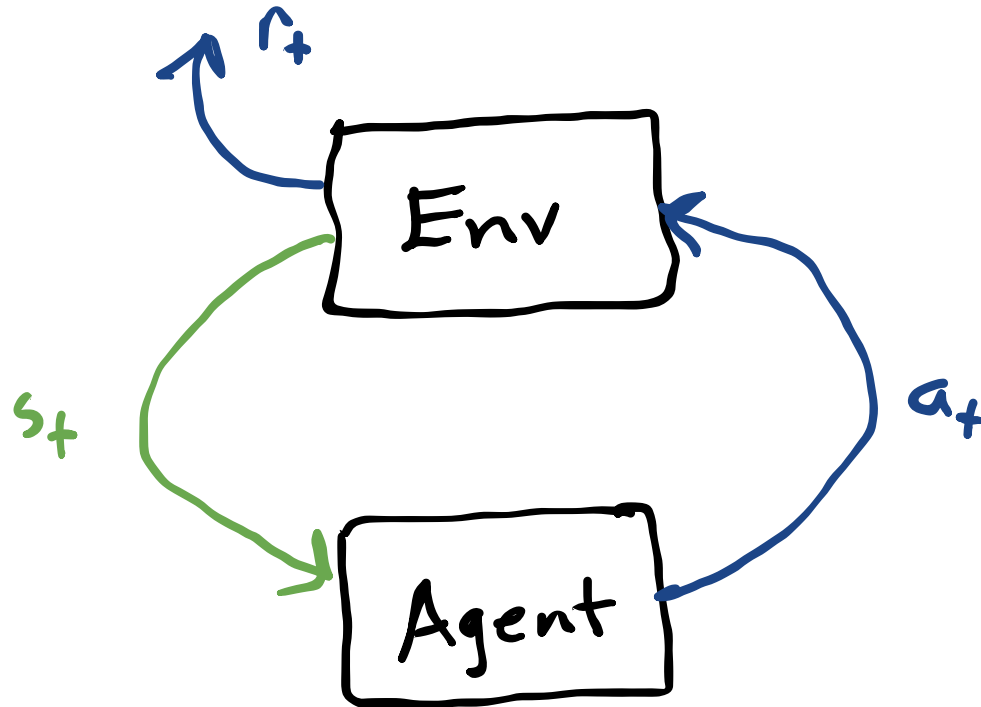
$r = \text{act!}(\text{env}, a)$

# Reinforcement Learning

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

Unknown!

Now: Episodic Simulator



`r = act!(env, a)`

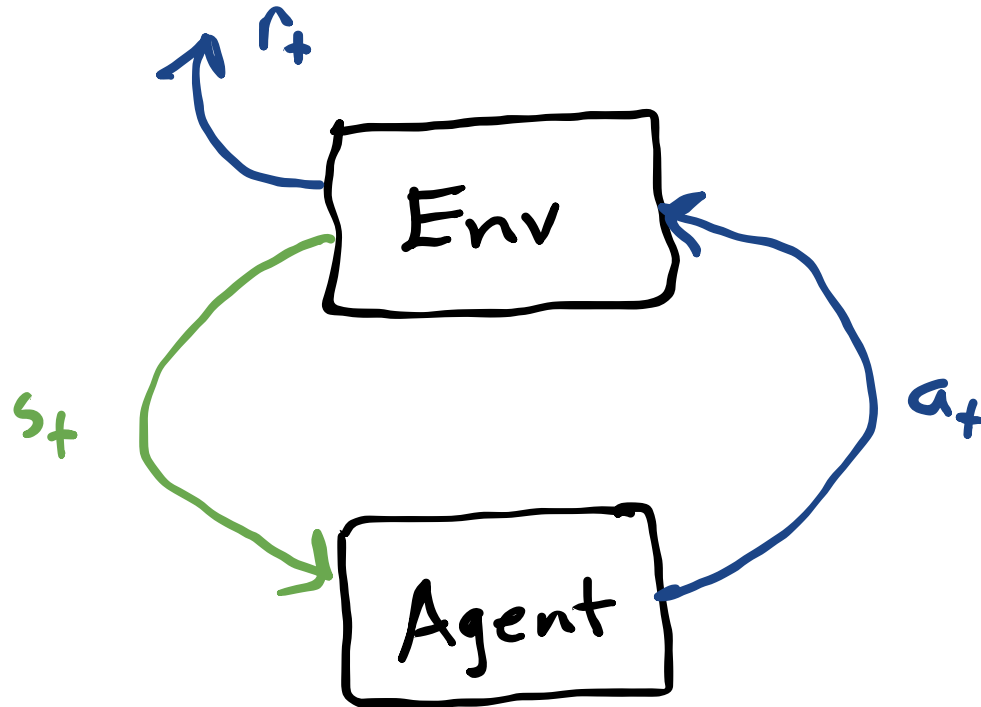
`s = observe(env)`



# Reinforcement Learning

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

Now: Episodic Simulator



```
r = act!(env, a)
```

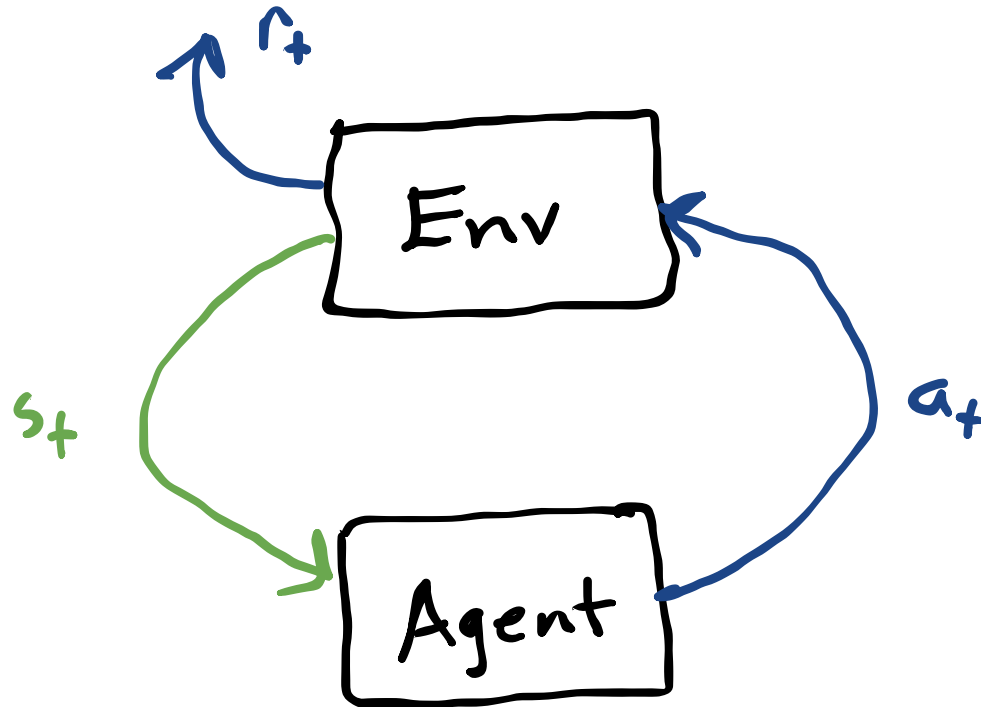
```
s = observe(env)
```

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)`

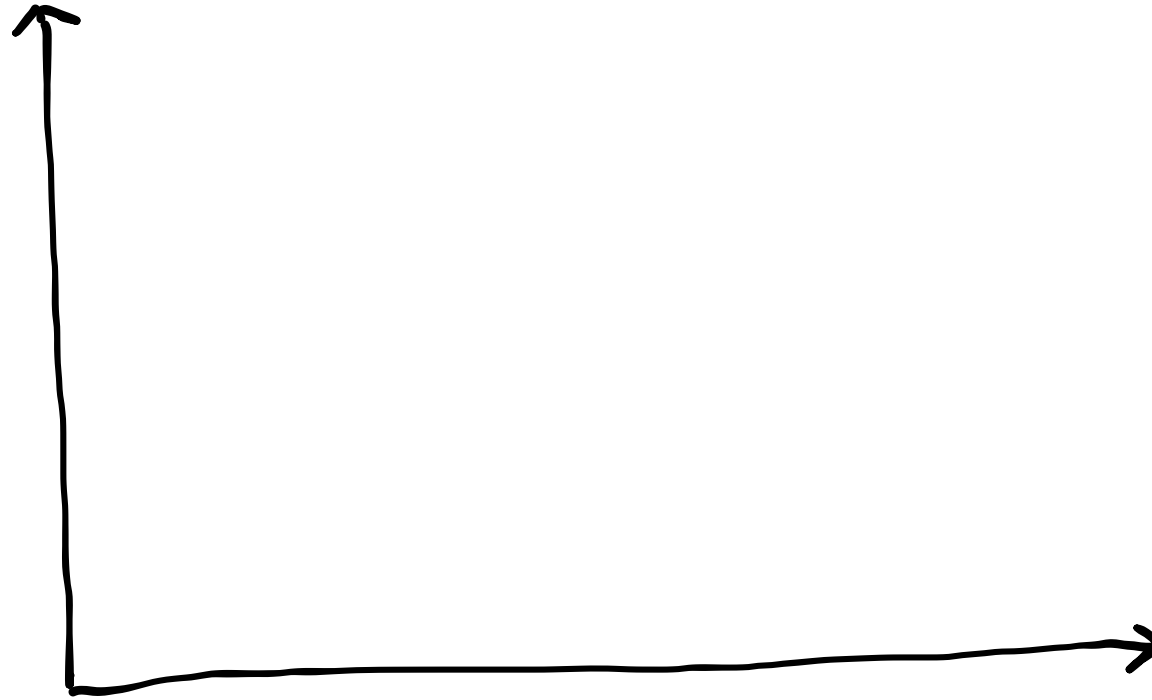
`s = observe(env)`

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

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

# Learning Curve

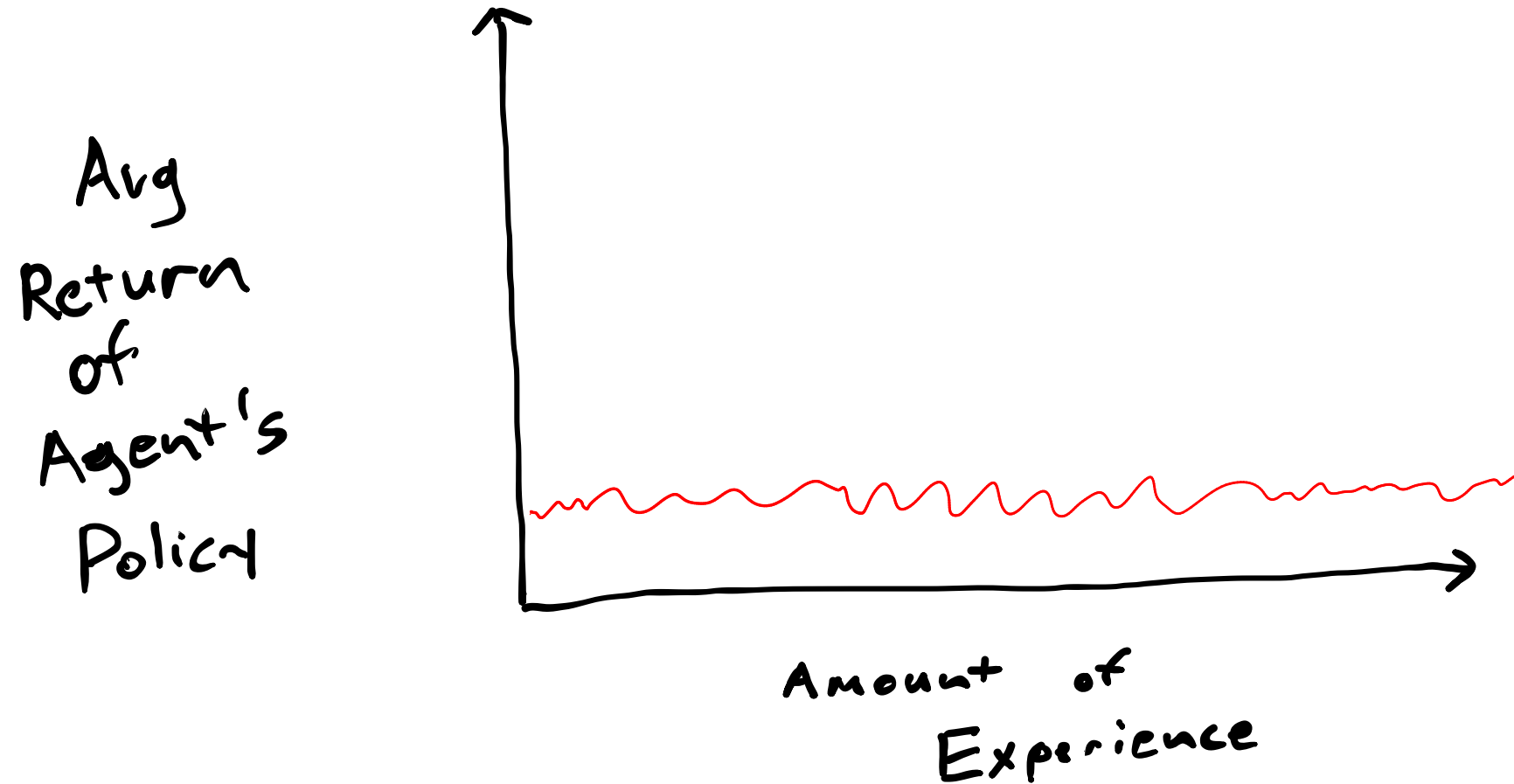
# Learning Curve



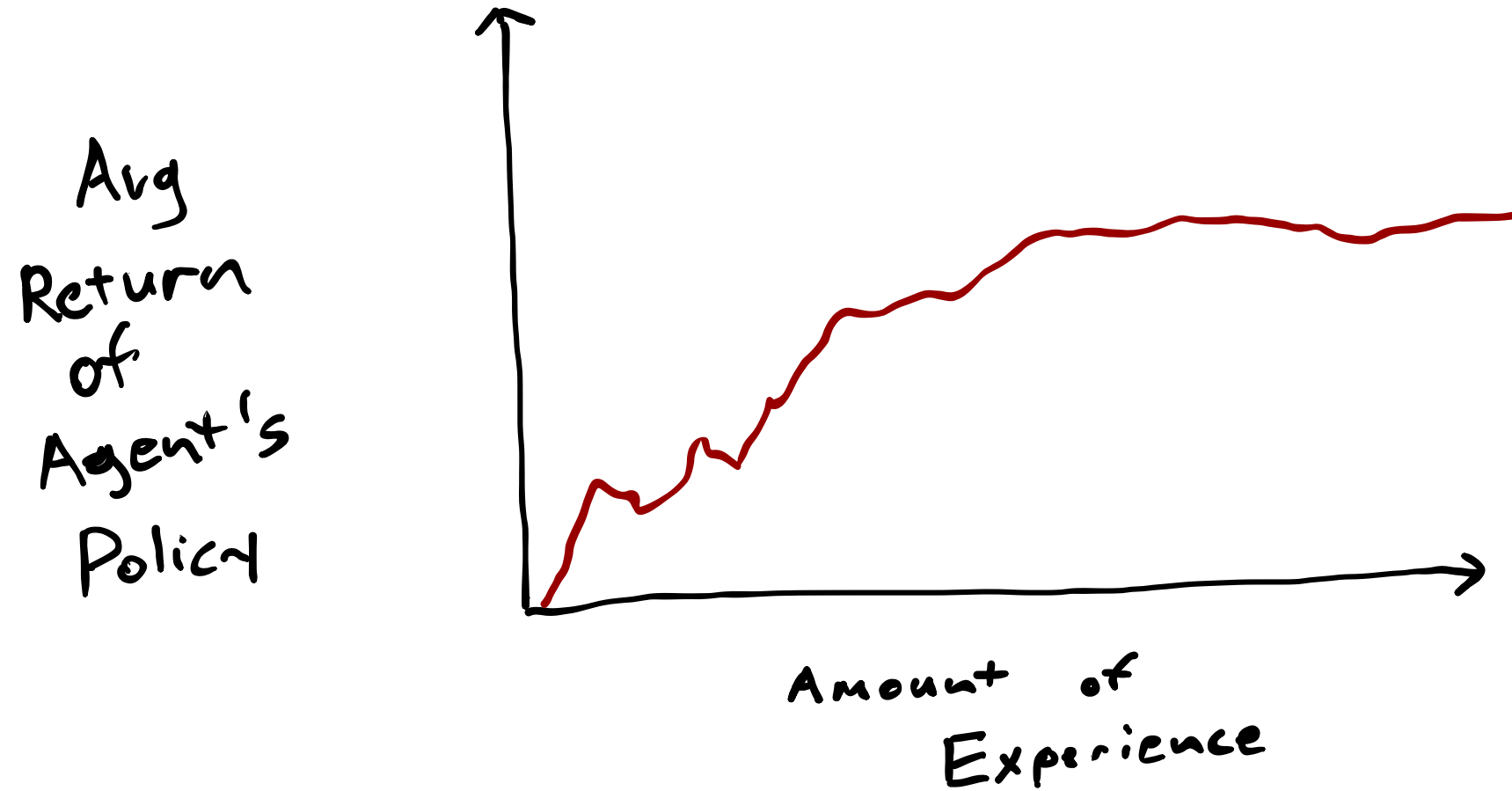
# Learning Curve



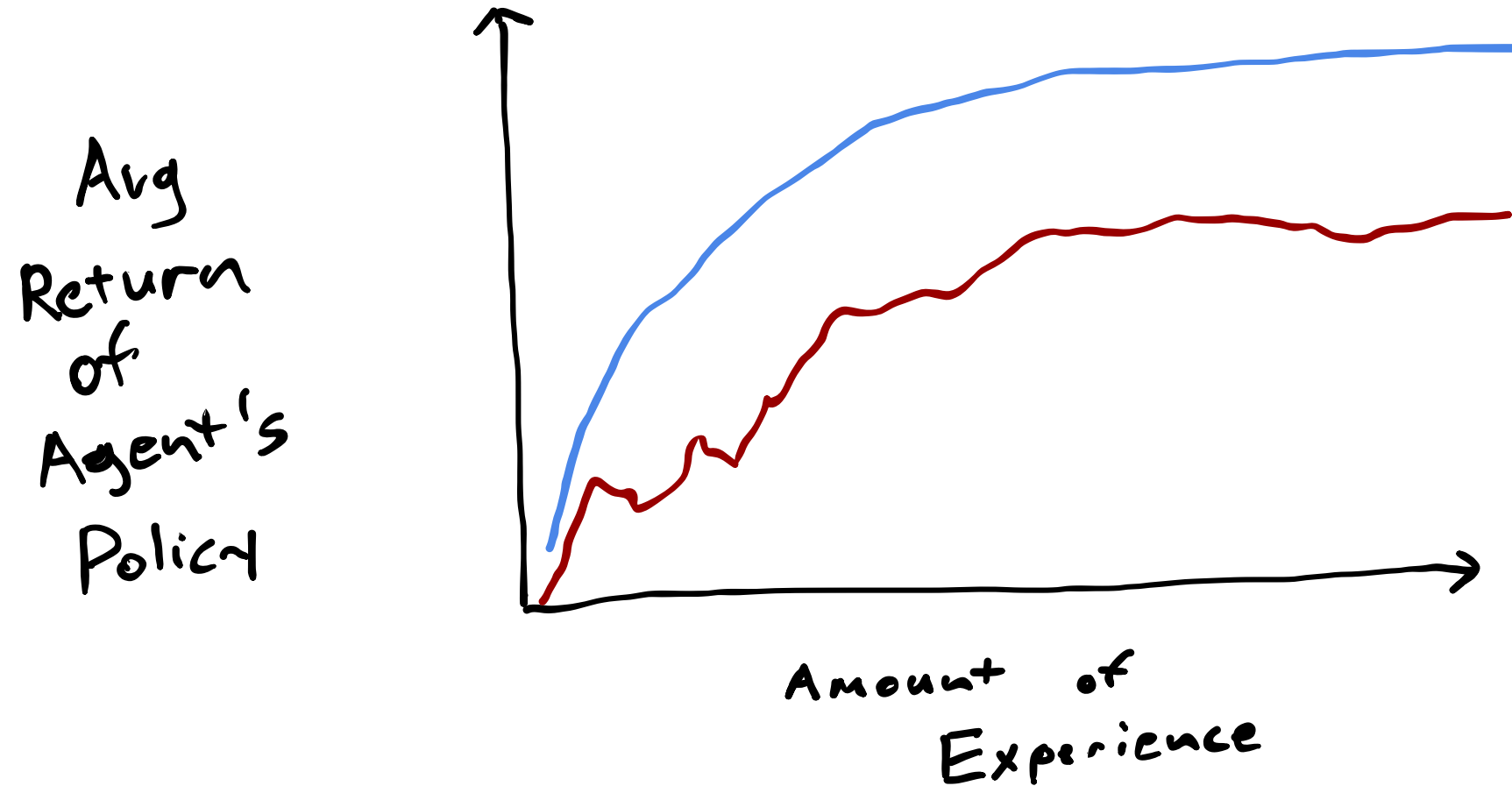
# Learning Curve



# Learning Curve

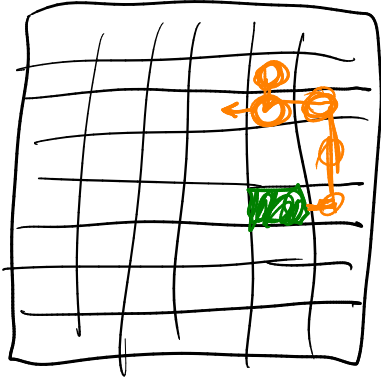


# Learning Curve





# Break



know  $S, A$

How should we interact with environment  
to maximize return given no knowledge of  $T$  and  $R$ .

model-based {  $\hat{R}(s, a)$   
 $\hat{N}(s, a, s')$

$$\longrightarrow \frac{1}{T}$$

Model-free

Actor Critic

$\hat{U}(s)$   
 $\hat{Q}(s, a)$   
 $\pi(s)$

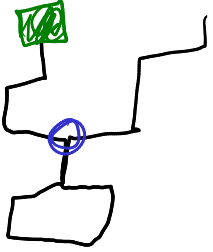
$$\pi(s) = \operatorname{argmax}_a \left( \hat{R}(s, a) + \gamma \underset{\hat{T}}{\mathbb{E}} [U(s)] \right)$$

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

# 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  $\pi^*$  by solving MDP

# Classifications

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



# Classifications

- **Model Based:** Attempt to learn  $T$  and  $R$ , then find  $\pi^*$  by solving MDP
- **Model Free:** Attempt to find  $Q^*$  or  $\pi^*$  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  $\pi^*$  by solving MDP
- **Model Free:** Attempt to find  $Q^*$  or  $\pi^*$  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  $\pi^*$  by solving MDP
- **Model Free:** Attempt to find  $Q^*$  or  $\pi^*$  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  $\pi^*$  by solving MDP
- **Model Free:** Attempt to find  $Q^*$  or  $\pi^*$  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  $\pi^*$  by solving MDP
- **Model Free:** Attempt to find  $Q^*$  or  $\pi^*$  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

TML MBRL

Given env,  $S, A$

$$\begin{aligned} N[s, a, s'] &\leftarrow 0 & \forall s, a, s' \\ \rho[s, a] &\leftarrow 0 & \forall s, a \end{aligned} \quad \left. \vphantom{\begin{aligned} N[s, a, s'] &\leftarrow 0 \\ \rho[s, a] &\leftarrow 0 \end{aligned}} \right\}$$

cumulative reward

$s \leftarrow \text{observe}(\text{env})$   
 $\pi \leftarrow \text{random policy}$

loop

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

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

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

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

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

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

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

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

$$s \leftarrow s'$$

$$\left. \begin{aligned} &\forall s, a, s' \\ &\forall s, a \end{aligned} \right\}$$

expensive

# Guiding Questions

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