

# Q-Learning

CMPT 729 G100

Jason Peng

# Overview

---

- Q-Function
- Q-Learning
- Exploration

# Taxonomy of RL Algorithms

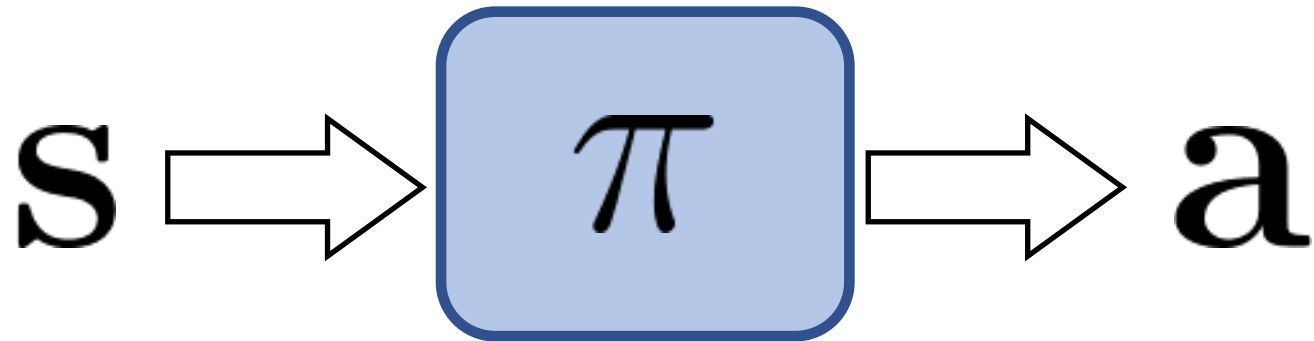
---

- Policy-Based Methods
- Value-Based Methods
- Actor-Critic Methods
- Model-Based Methods

# Policy-Based Methods

---

$$\pi(\mathbf{a}|\mathbf{s})$$



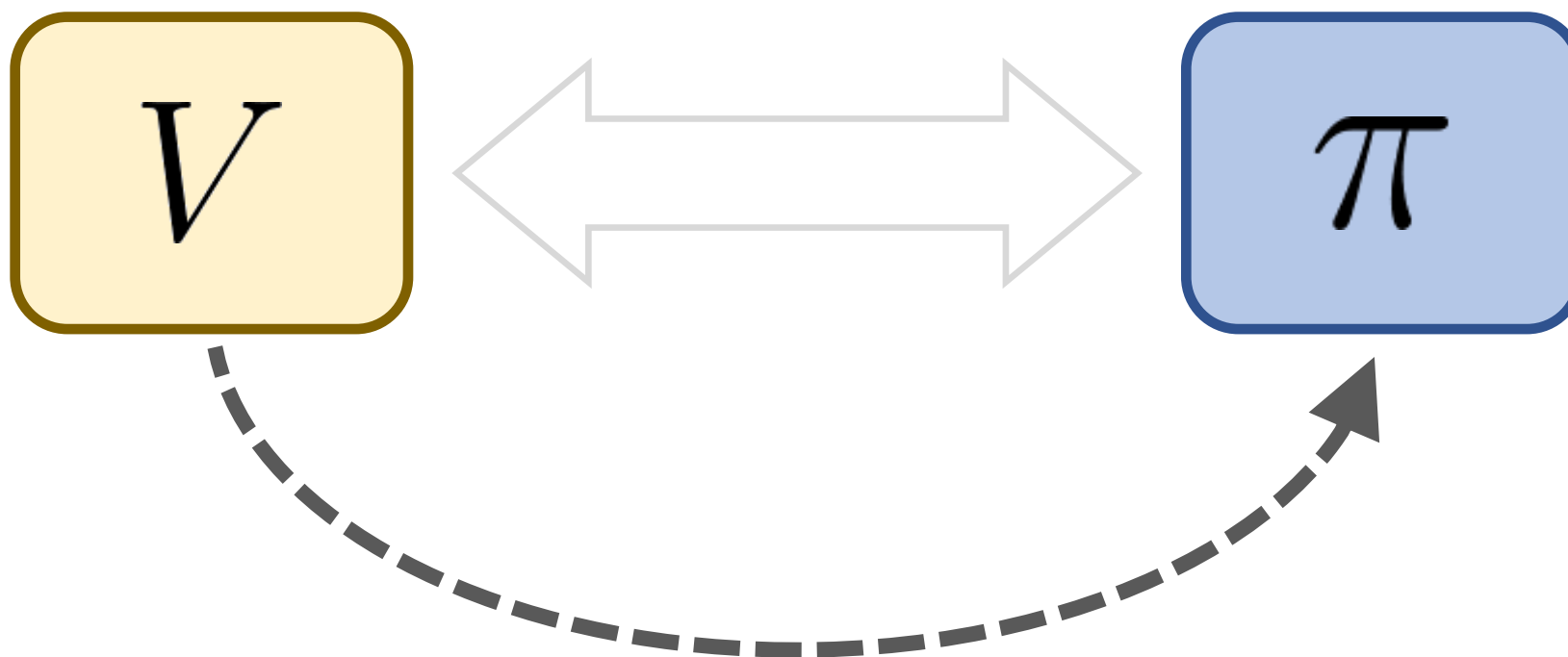
# Value-Based Methods

---



# Value-Based Methods

---




# Policy Gradient

---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\tau \sim p(\tau|\pi, \mathbf{s}_0=\mathbf{s}, \mathbf{a}_0=\mathbf{a})} \left[ \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \left( \sum_{t=0}^{\tau} \gamma^t r_t - \cancel{V^{\pi}(\mathbf{s})} \right) \right]$$

# Policy Gradient

---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\tau \sim p(\tau|\pi, \mathbf{s}_0=\mathbf{s}, \mathbf{a}_0=\mathbf{a})} \left[ \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \left( \sum_{t=0}^{\tau} \gamma^t r_t \right) \right]$$




# Policy Gradient

---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[ \underbrace{\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \mathbb{E}_{\tau \sim p(\tau|\pi, \mathbf{s}_0=\mathbf{s}, \mathbf{a}_0=\mathbf{a})} \left( \sum_{t=0}^{\infty} \gamma^t r_t \right)}_{\substack{\text{“reward-to-go”} \\ = Q^{\pi}(\mathbf{s}, \mathbf{a})}} \right]$$

reward-to-go: expected return of taking an action  $\mathbf{a}$  in state  $\mathbf{s}$

# Policy Gradient

---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{Q^{\pi}(\mathbf{s}, \mathbf{a})}]$$

reward-to-go: expected return of taking an action  $\mathbf{a}$  in state  $\mathbf{s}$

# Policy Gradient

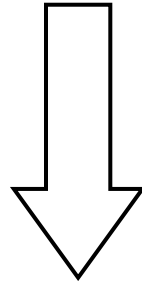
---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\underline{\mathbf{s}} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) Q^{\pi}(\mathbf{s}, \mathbf{a})]$$

# Policy Gradient

---

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) Q^{\pi}(\mathbf{s}, \mathbf{a})]$$

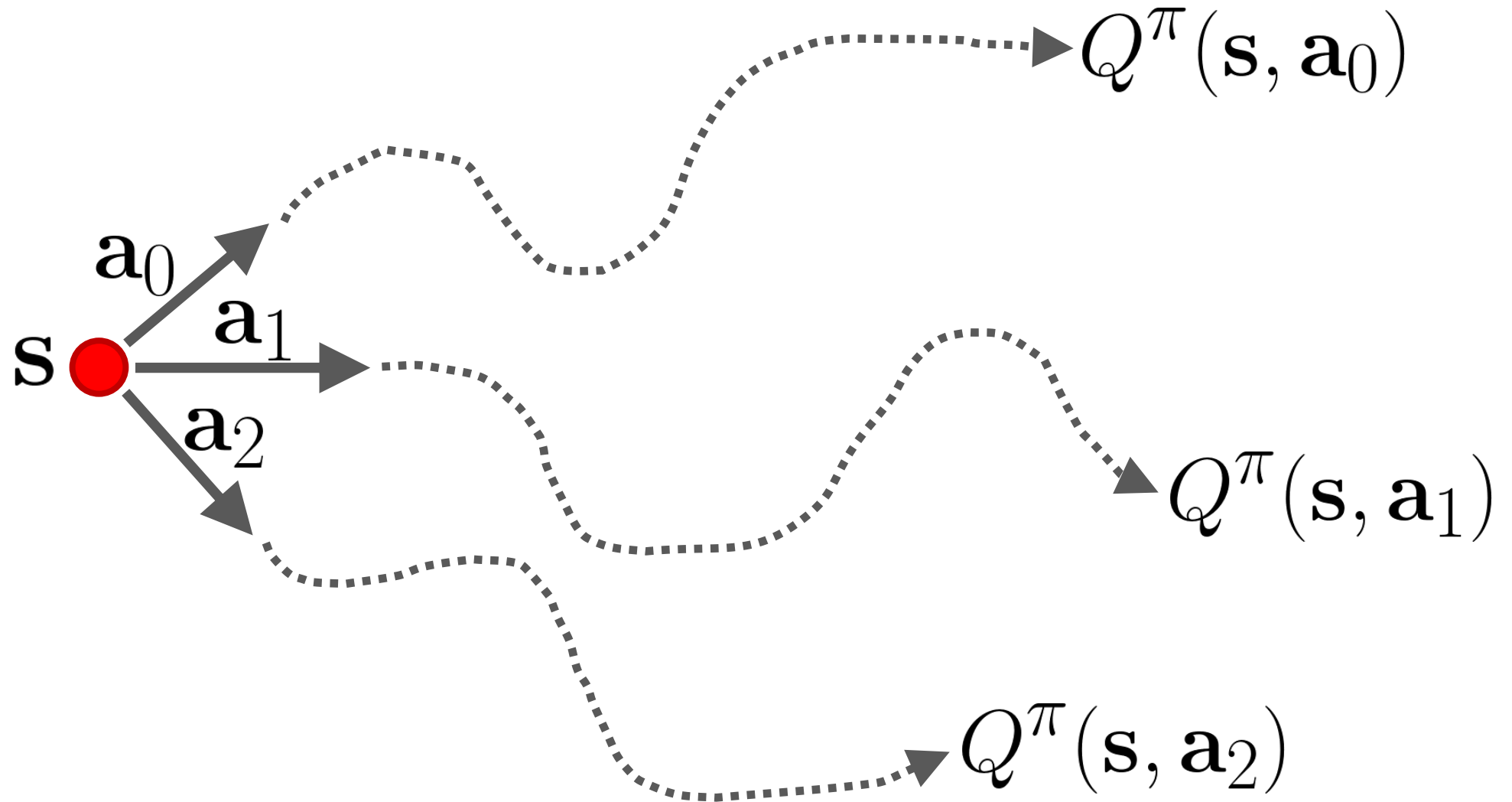


$$\max_{\pi} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^{\pi}(\mathbf{s}, \mathbf{a})]$$

**Per-state objective:** pick actions that maximize the expected return at each state (i.e. Q-function)

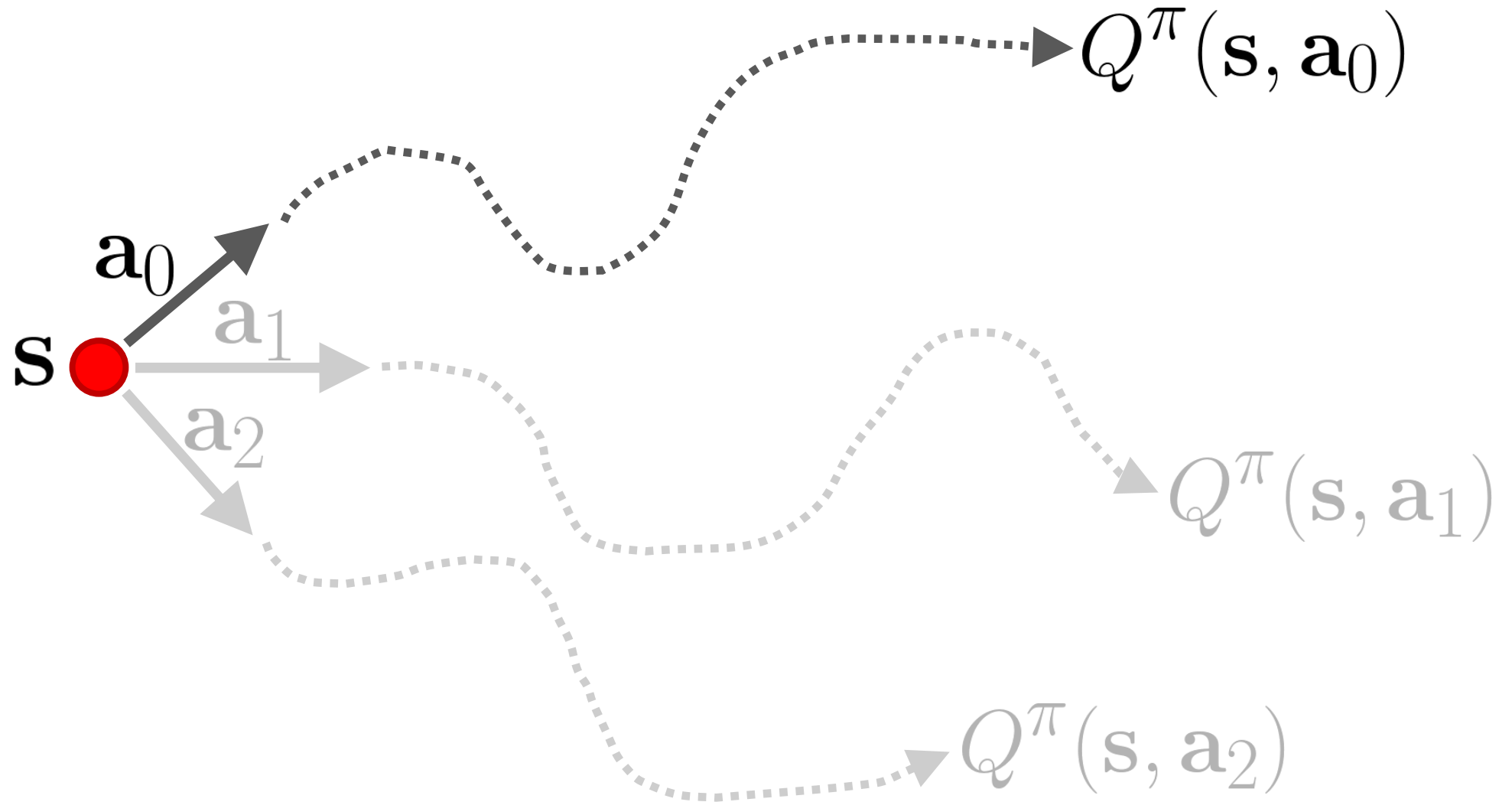
# Q-Function

---



# Q-Function

---




# Value Functions

---

## Value Function


“State Value Function”

$$V^{\pi}(\mathbf{s}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$


Likelihood of a trajectory  
starting at state **S** and  
then following  $\pi$  for all  
future timesteps

## Q-Function

“State-Action Value Function”

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$


Likelihood of a trajectory  
after taking action **a** in  
state **S** and then following  
 $\pi$  for all future timesteps

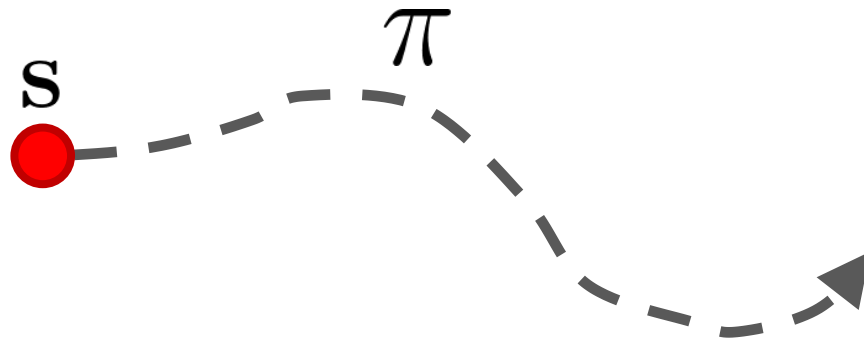
# Value Functions

---

## Value Function

“State Value Function”

$$V^{\pi}(\mathbf{s}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

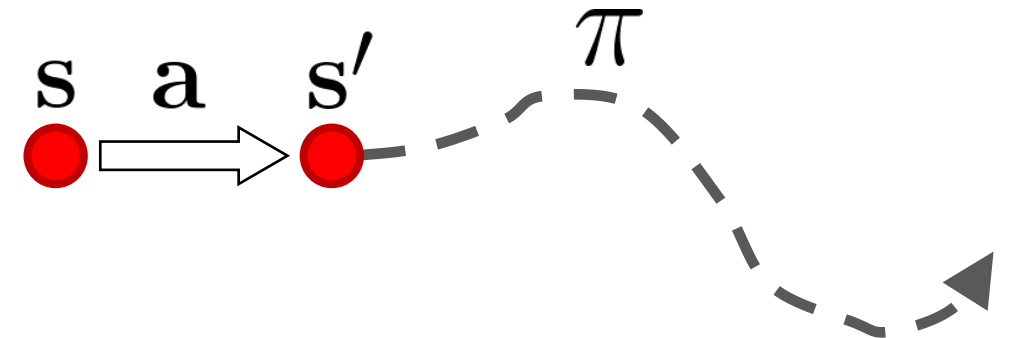


## Q-Function

“State-Action Value Function”

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

“Quality”





# Recursive Definition

---

## Value Function

$$V^{\pi}(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [\underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma V^{\pi}(\mathbf{s}')] ]$$

# Recursive Definition

---

## Value Function

$$V^{\pi}(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \underline{V^{\pi}(\mathbf{s}')}]$$

# Recursive Definition

---

## Value Function

$$V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] ]$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [\underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] ]$$

# Recursive Definition

---

## Value Function

$$V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] ]$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [\underline{Q^\pi(\mathbf{s}', \mathbf{a}')}]] \right]$$

# Recursive Definition

---

## Value Function

$$V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] ]$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

# Recursive Definition

---

## Value Function

$$V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] ]$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

# Recursive Definition

---

## Value Function

$$V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \underbrace{\mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] }_{= Q^\pi(\mathbf{s}, \mathbf{a})}$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

# Recursive Definition

---

## Value Function

$$\begin{aligned} V^\pi(\mathbf{s}) &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})] \end{aligned}$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$



# Recursive Definition

---

## Value Function

$$\begin{aligned} V^\pi(\mathbf{s}) &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})] \end{aligned}$$

## Q-Function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [\underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] ]$$

# Recursive Definition

---

## Value Function

$$\begin{aligned} V^\pi(\mathbf{s}) &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})] \end{aligned}$$

## Q-Function

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \underbrace{\gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] }_{= V^\pi(\mathbf{s}')} \right] \end{aligned}$$

# Recursive Definition

---

## Value Function

$$\begin{aligned} V^\pi(\mathbf{s}) &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] \\ &= \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})] \end{aligned}$$

## Q-Function

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right] \\ &= \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^\pi(\mathbf{s}')] \end{aligned}$$

# Q-Function

---

$$\cancel{\pi^*(\mathbf{a}|\mathbf{s})} \Rightarrow Q^*(\mathbf{s}, \mathbf{a})$$

Recover optimal policy:

$$\pi^*(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^*(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$

Instead of learning policy, just learn Q-function.

# Q-Function

---

$$\pi(\mathbf{a}|\mathbf{s}) \Longrightarrow Q^{\pi}(\mathbf{s}, \mathbf{a})$$

Recover a policy:

“arg max policy”

$$\pi'(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^{\pi}(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$

New policy is at least as good as the old policy.

$$J(\pi') \geq J(\pi) \quad Q^{\pi'}(\mathbf{s}, \mathbf{a}) \geq Q^{\pi}(\mathbf{s}, \mathbf{a})$$

# Q-Learning

---

Key idea:

- Instead of trying to learn the optimal policy, just learn optimal Q-function
- Then recover policy from Q-function

# Q-Learning

---

## Recursive definition

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

## Optimal policy

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi^*(\mathbf{a}'|\mathbf{s}')} [Q^*(\mathbf{s}', \mathbf{a}')] \right]$$

$$\pi^*(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^*(\mathbf{s}, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$$

# Q-Learning

---

Recursive definition

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

Optimal policy

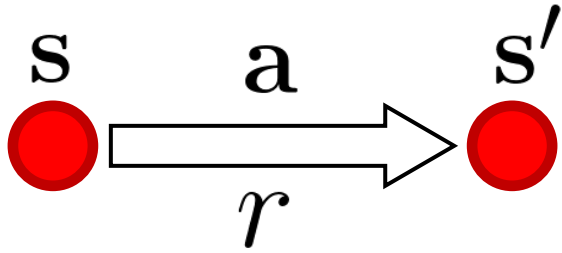
$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$



# Q-Learning

---

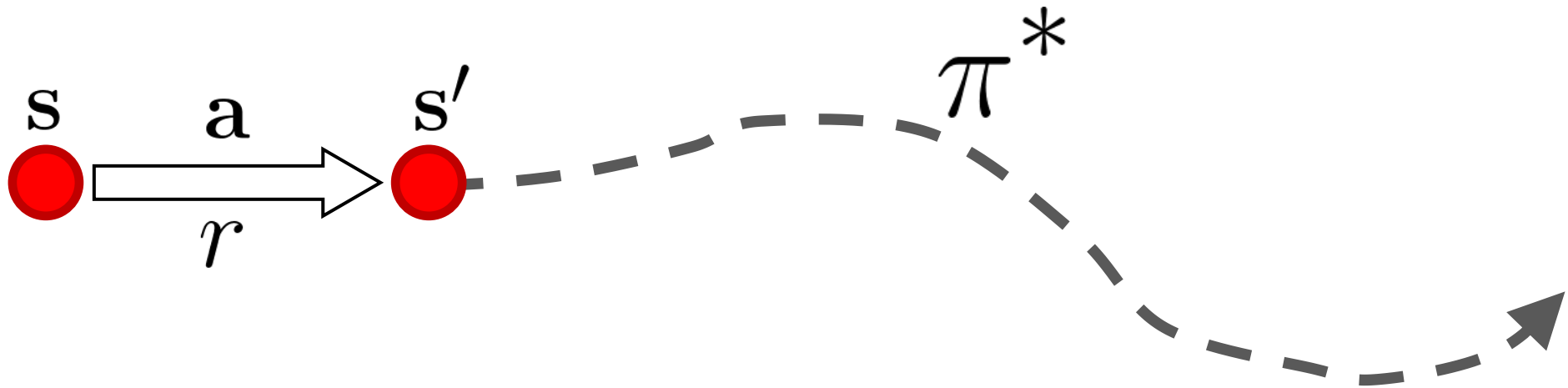
$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$



# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$



# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \neq \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]$$

# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]$$

$\geq$

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq \underbrace{\mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]}$$

# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]$$

# Q-Learning

---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]$$

arg max policy

# Q-Learning

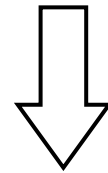
---

$$Q^*(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^*(\mathbf{s}', \mathbf{a}')) \right]$$

Not true for non-optimal policies

$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq \underbrace{\mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} (Q^\pi(\mathbf{s}', \mathbf{a}')) \right]}_{Q^{\pi'}(\mathbf{s}, \mathbf{a})}$$

$$Q^{\pi'}(\mathbf{s}, \mathbf{a})$$



$$Q^\pi(\mathbf{s}, \mathbf{a}) \leq Q^{\pi'}(\mathbf{s}, \mathbf{a})$$



# Q-Learning

---

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\underline{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})}} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

# Q-Learning

---

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

# Q-Learning

---

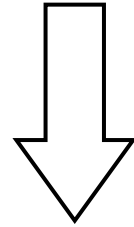
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

# Q-Learning

---

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) \geq Q^k(\mathbf{s}, \mathbf{a})$$









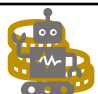
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = Q^k(\mathbf{s}, \mathbf{a})$$

$$Q^k(\mathbf{s}, \mathbf{a}) = Q^*(\mathbf{s}, \mathbf{a})$$

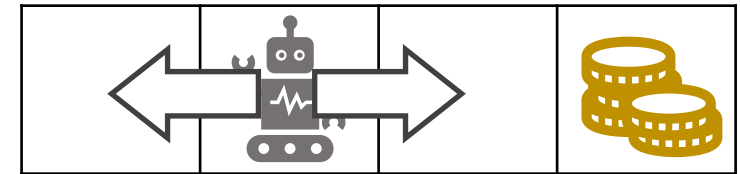
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment







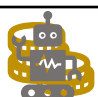


$$\gamma = 1/2$$

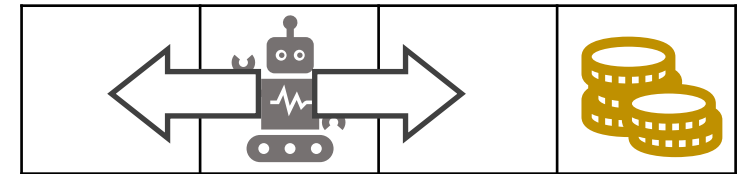
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment










$$\gamma = 1/2$$

# Tabular Q-Learning

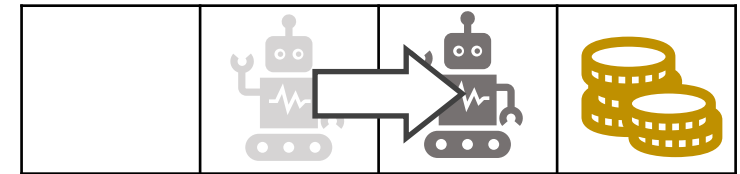
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

$= 0$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment










$$\gamma = 1/2$$

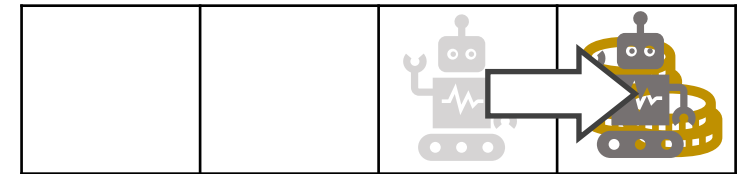
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment



$$\gamma = 1/2$$











# Tabular Q-Learning

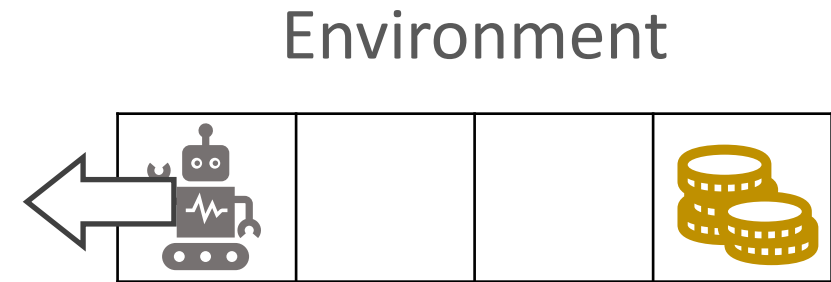
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

$= 0$

Iteration 0:

$Q$

				Action		
				←	□	→
State					0	0
					0	0
					0	0
					0	0



$$\gamma = 1/2$$








# Tabular Q-Learning

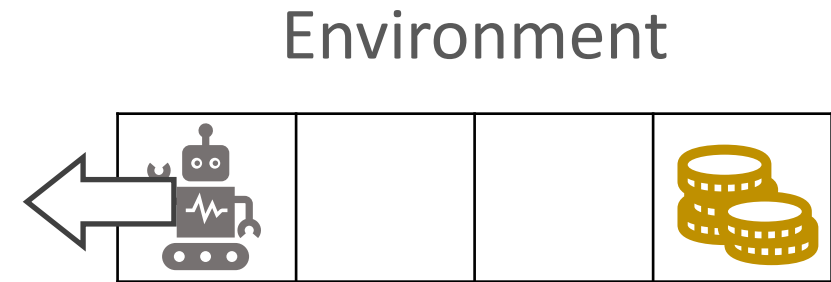
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

$Q$

State

				Action		
				←	□	→
				0	0	0
				0	0	0
				0	0	0
				0	0	0



$$\gamma = 1/2$$








# Tabular Q-Learning

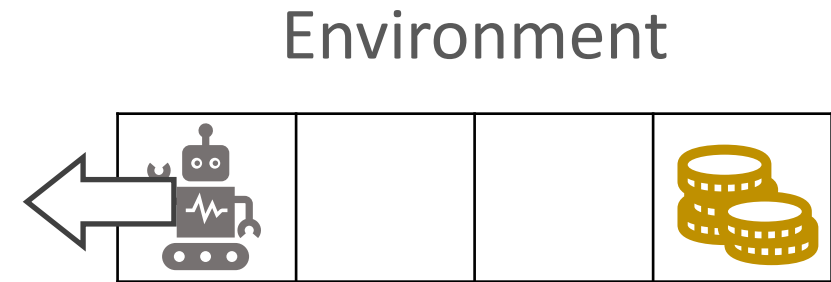
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State

$Q$

				Action		
				←	□	→
				0	0	0
				0	0	0
				0	0	0
				0	0	0



$$\gamma = 1/2$$

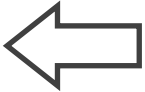

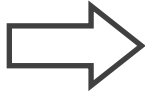










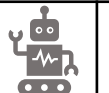
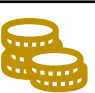
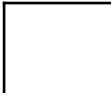
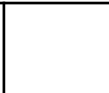

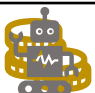
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')}_{= 0} \right) \right]$$

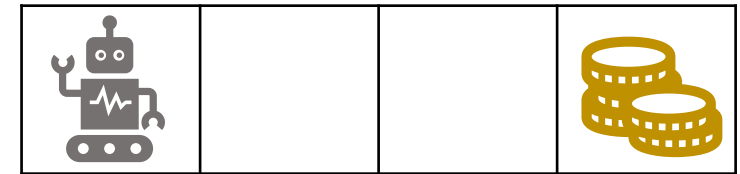
Iteration 0:

State

$Q$

Action			
<div>    </div>			
   	0	0	0
   	0	0	0
   	0	0	0
   	0	0	0

Environment

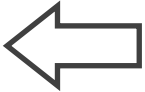

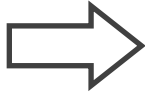


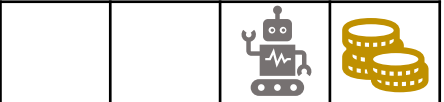



$$\gamma = 1/2$$

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

		Action		
				
State		0	0	0
		0	0	0
		0	0	0
		0	0	0

Environment



$$\gamma = 1/2$$

# Tabular Q-Learning

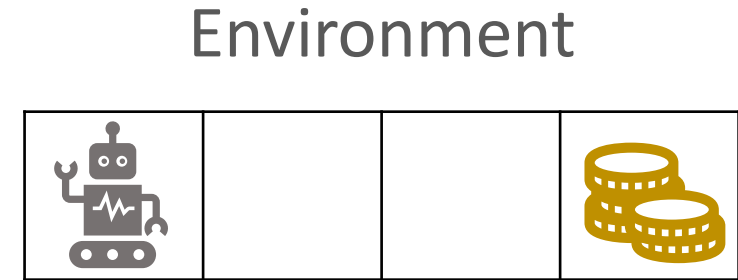
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State

$Q$

Action	←	□	→
Robot at start, coins at end	0	0	0
Robot in middle, coins at end	0	0	0
Robot at end, coins at end	0	0	0
Robot at end, coins at start	0	0	0



$$\gamma = 1/2$$








# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underline{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

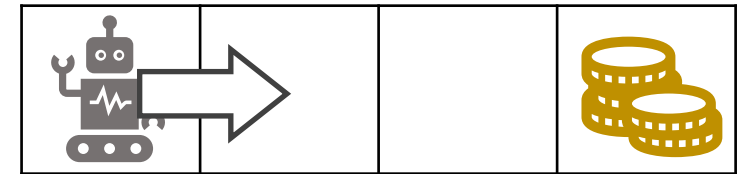
Iteration 0:

State

$Q$

				Action		
				←	□	→
				0	0	0
				0	0	0
				0	0	0
				0	0	0

Environment



$$\gamma = 1/2$$




# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')}\right) \right]$$

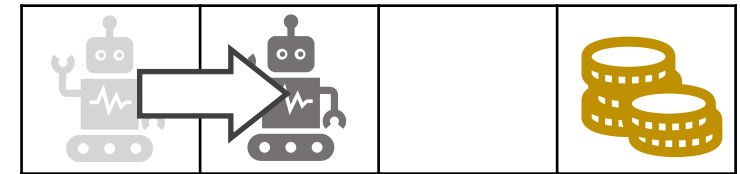
$= 0$

Iteration 0:

$Q$

				Action		
				←	□	→
State				0	0	0
				0	0	0
				0	0	0
				0	0	0

Environment



$$\gamma = 1/2$$









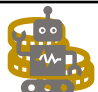
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

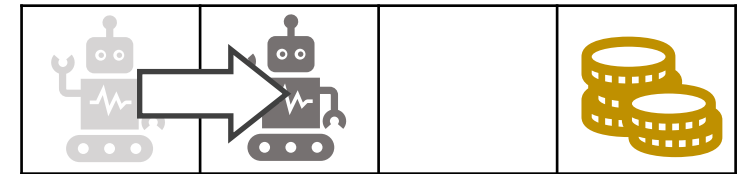
Iteration 0:

State

$Q$

				Action		
				←	□	→
				0	0	0
				0	0	0
				0	0	0
				0	0	0

Environment



$$\gamma = 1/2$$

# Tabular Q-Learning

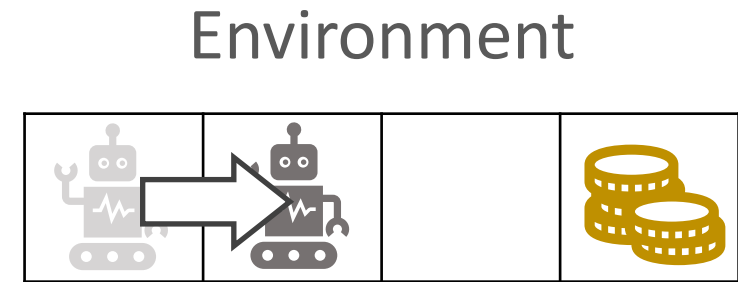
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State

$Q$

Action	←	□	→
Robot at left, coins at right	0	0	0
Robot at second from left, coins at right	0	0	0
Robot at third from left, coins at right	0	0	0
Robot at right, coins at right	0	0	0



$$\gamma = 1/2$$










# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

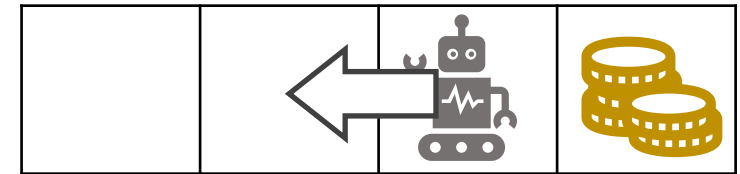
Iteration 0:

State

$Q$

				Action		
				←	□	→
				0	0	0
				0	0	0
				0	0	0
				0	0	0

Environment







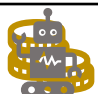


$$\gamma = 1/2$$

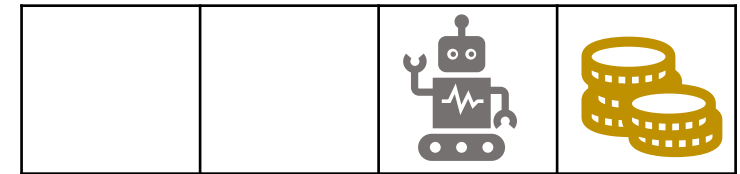
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment


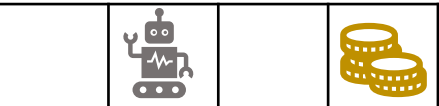
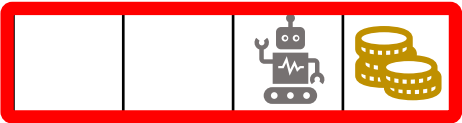



$$\gamma = 1/2$$

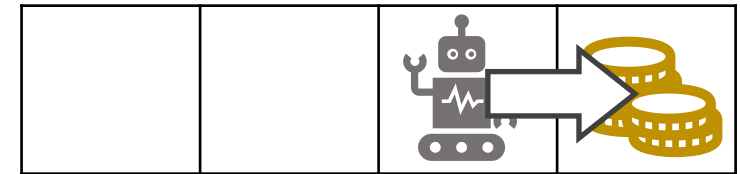
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	0
					0	0	0

Environment



$$\gamma = 1/2$$


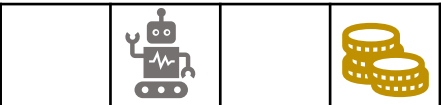
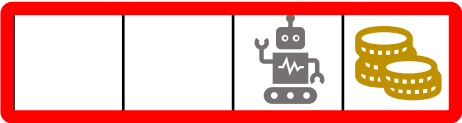

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

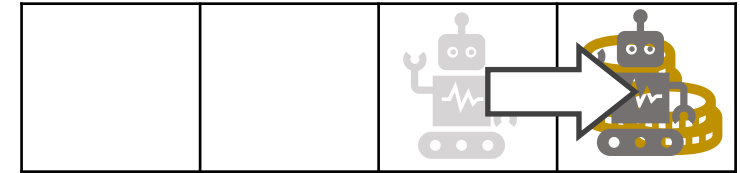
Iteration 0:

State

$Q$

	Action		
	←	□	→
	0	0	0
	0	0	0
	0	0	0
	0	0	0

Environment



$$\gamma = 1/2$$


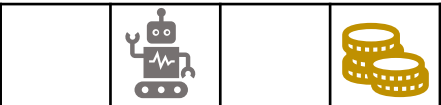
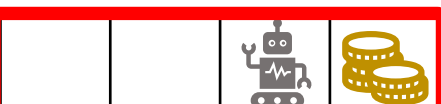
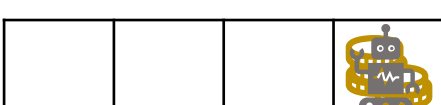
# Tabular Q-Learning

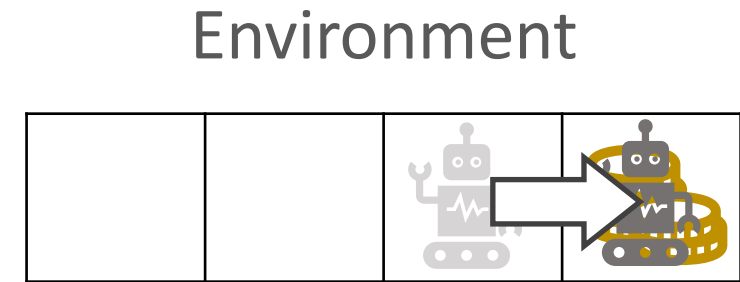
$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State

$Q$

	Action		
	←	□	→
	0	0	0
	0	0	0
	0	0	0
	0	0	0



$$\gamma = 1/2$$




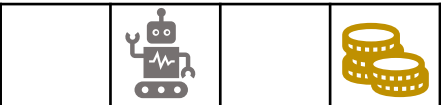
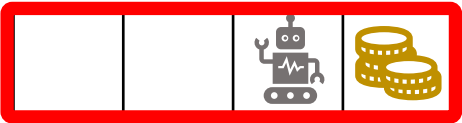

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} (Q^k(\mathbf{s}', \mathbf{a}'))}_{= 0} \right]$$

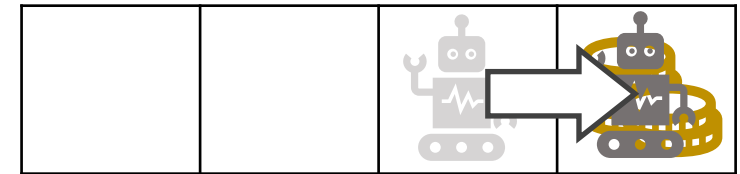
Iteration 0:

State

$Q$

	Action		
	←	□	→
	0	0	0
	0	0	0
	0	0	1
	0	0	0

Environment



$$\gamma = 1/2$$

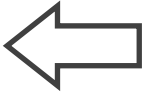

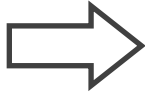















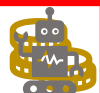
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')} \right) \right]$$

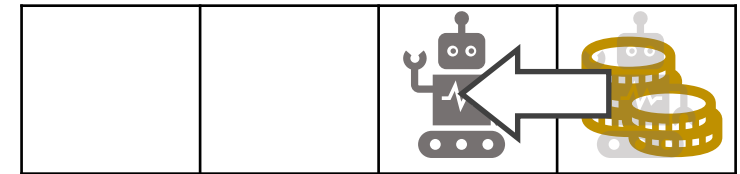
$= 0$

Iteration 0:

$Q$

State	Action			
				
	   	0	0	0
	   	0	0	0
	   	0	0	1
	   	0	0	0

Environment



$$\gamma = 1/2$$


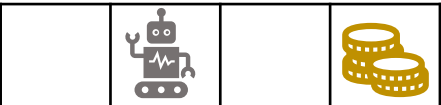
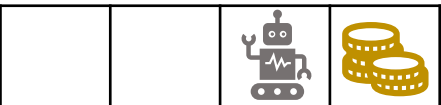

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \underbrace{\max_{\mathbf{a}'} (Q^k(\mathbf{s}', \mathbf{a}'))}_{= 1} \right]$$

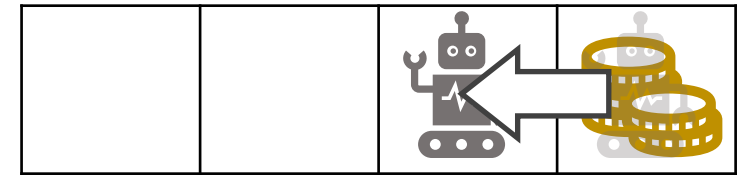
Iteration 0:

State

$Q$

	Action		
	←	□	→
	0	0	0
	0	0	0
	0	0	1
	0	0	0

Environment


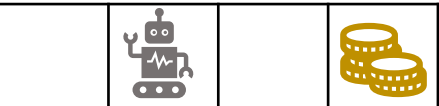
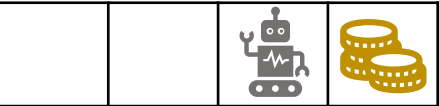



$$\gamma = 1/2$$

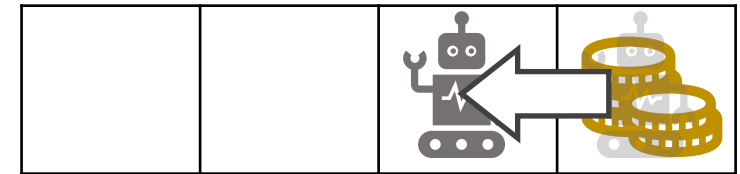
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 1} \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	0	0

Environment


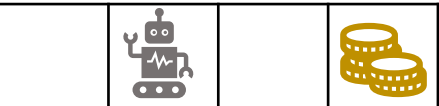
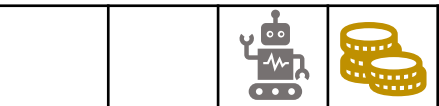
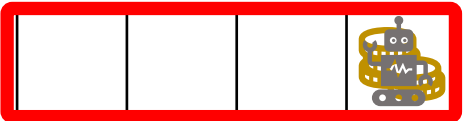


$$\gamma = 1/2$$

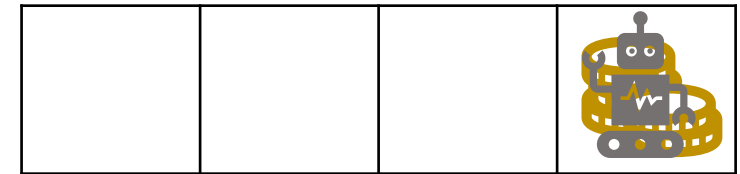
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	0	0

Environment


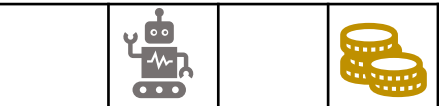
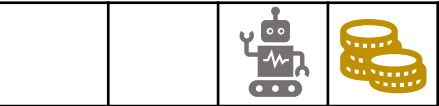
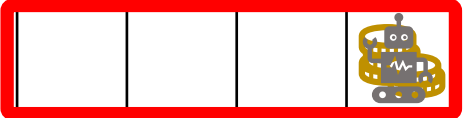


$$\gamma = 1/2$$

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
		1/2	0	0			

Environment


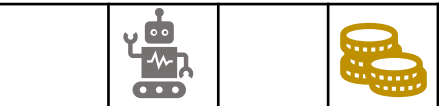
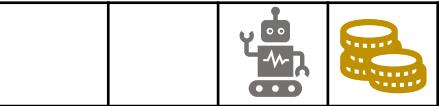
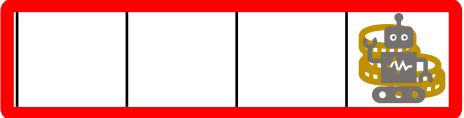


$$\gamma = 1/2$$

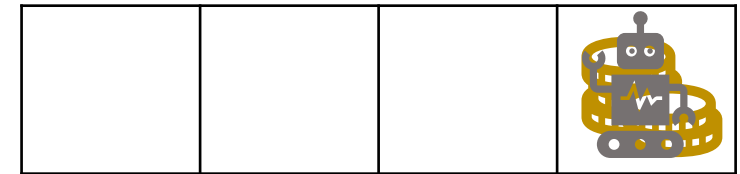
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} (Q^k(\mathbf{s}', \mathbf{a}'))}_{= 0} \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	0	0

Environment



$$\gamma = 1/2$$


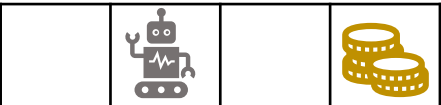
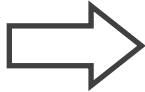

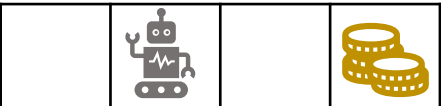
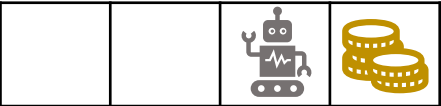
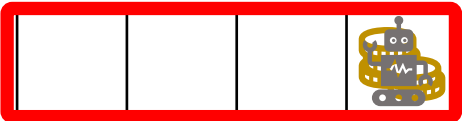
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

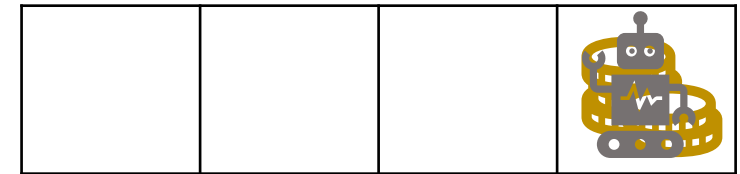
Iteration 0:

State

$Q$

	Action		
			
	0	0	0
	0	0	0
	0	0	1
	1/2	1	0

Environment




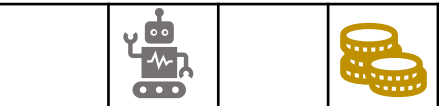
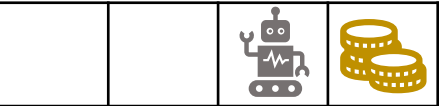

$$\gamma = 1/2$$



# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')} \right) \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	1	0

Environment


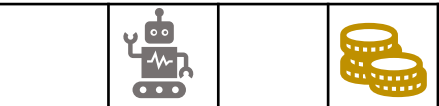
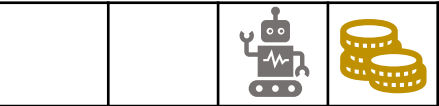



$$\gamma = 1/2$$

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	1	0

Environment


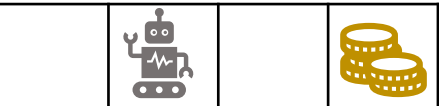
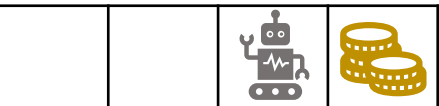



$$\gamma = 1/2$$

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 1} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 0} \right]$$

Iteration 0:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
		1/2	1	1			

Environment









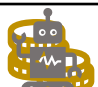
$$\gamma = 1/2$$



# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	1	1

Environment


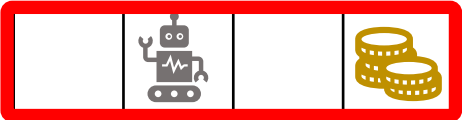
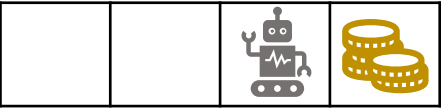



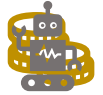


$$\gamma = 1/2$$

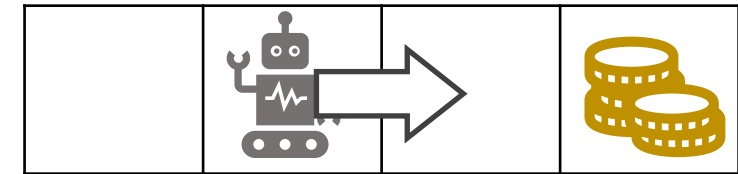
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	1	1

Environment



$$\gamma = 1/2$$


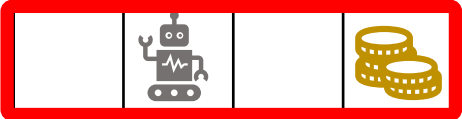
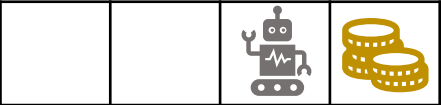

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')}_{= 0} \right) \right]$$

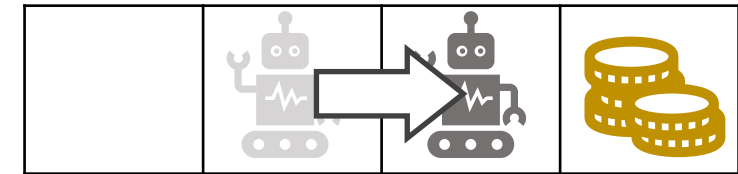
Iteration 1:

State

$Q$

	Action		
	←	□	→
	0	0	0
	0	0	0
	0	0	1
	1/2	1	1

Environment


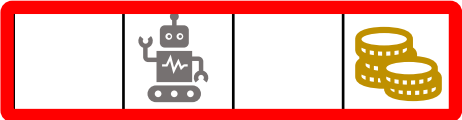
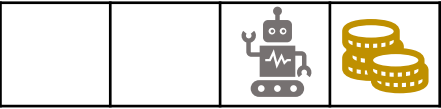



$$\gamma = 1/2$$

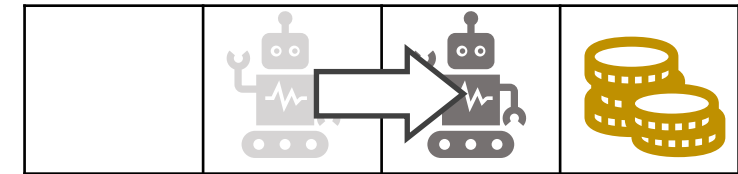
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 1} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	0
					0	0	1
					1/2	1	1

Environment




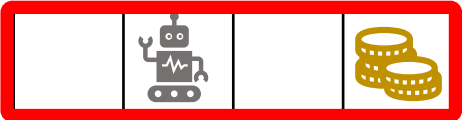
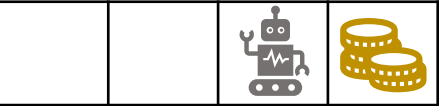

$$\gamma = 1/2$$



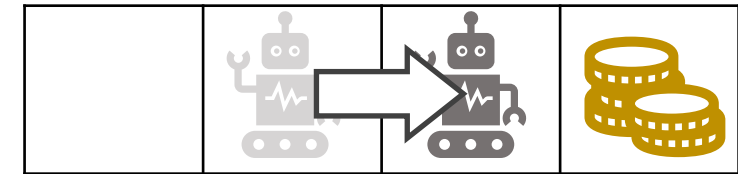
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \underbrace{\max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 1} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					0	0	1
					1/2	1	1

Environment


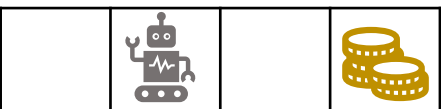
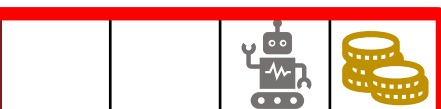



$$\gamma = 1/2$$

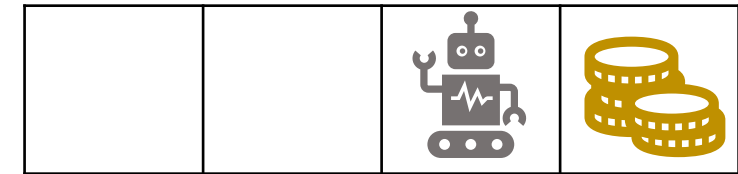
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					0	0	1
					1/2	1	1

Environment


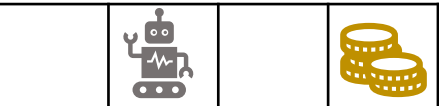
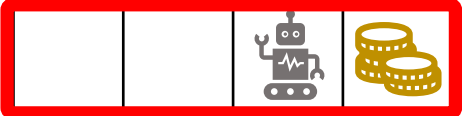



$$\gamma = 1/2$$

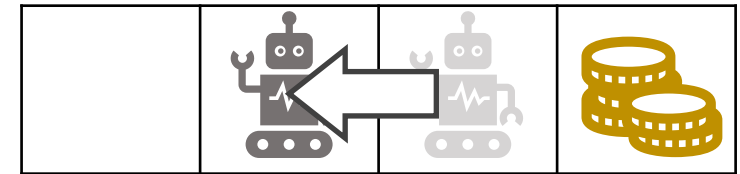
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \left( \underbrace{Q^k(\mathbf{s}', \mathbf{a}')}_{= 0} \right) \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					0	0	1
					1/2	1	1

Environment


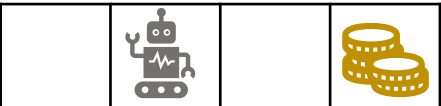
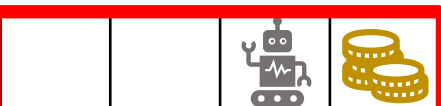



$$\gamma = 1/2$$

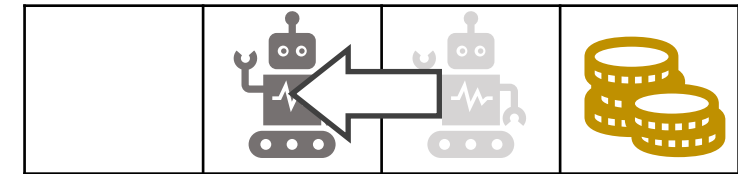
# Tabular Q-Learning

$$Q^{k+1}(s, a) = \mathbb{E}_{s' \sim p(s'|s, a)} \left[ \underbrace{r(s, a, s')}_{= 0} + \gamma \underbrace{\max_{a'} (Q^k(s', a'))}_{= 1/2} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					0	0	1
					1/2	1	1

Environment


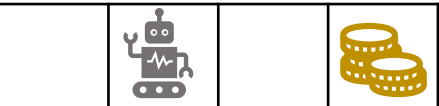
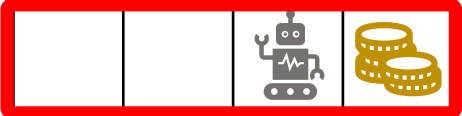



$$\gamma = 1/2$$

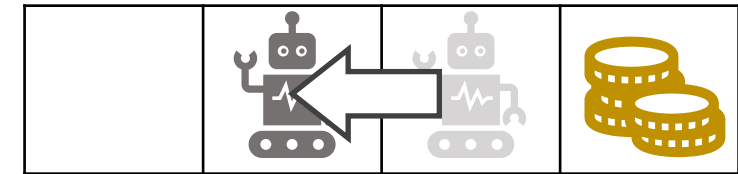
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \max_{\mathbf{a}'} \underbrace{\left( Q^k(\mathbf{s}', \mathbf{a}') \right)}_{= 1/2} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					1/4	0	1
					1/2	1	1

Environment


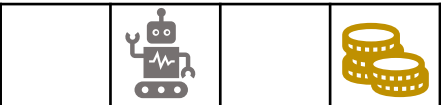
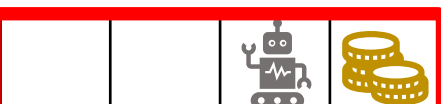
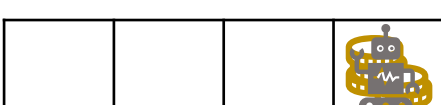


$$\gamma = 1/2$$

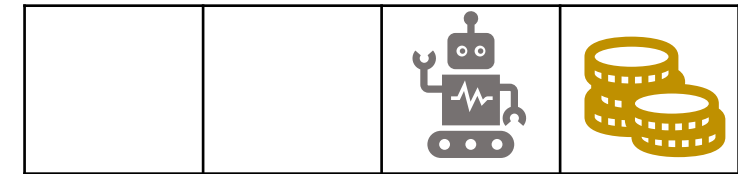
# Tabular Q-Learning

$$Q^{k+1}(s, a) = \mathbb{E}_{s' \sim p(s'|s, a)} \left[ \underbrace{r(s, a, s')}_{=0} + \gamma \max_{\underbrace{a'}} \left( Q^k(s', a') \right) \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
		1/4	0	1			
		1/2	1	1			

Environment


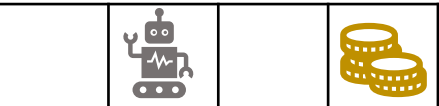
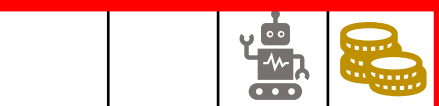



$$\gamma = 1/2$$

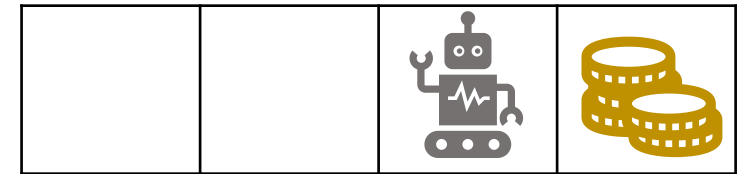
# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ \underbrace{r(\mathbf{s}, \mathbf{a}, \mathbf{s}')}_{= 0} + \gamma \underbrace{\max_{\mathbf{a}'} (Q^k(\mathbf{s}', \mathbf{a}'))}_{= 1} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
					1/4	0	1
					1/2	1	1

Environment


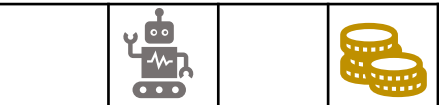
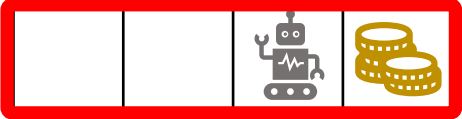



$$\gamma = 1/2$$

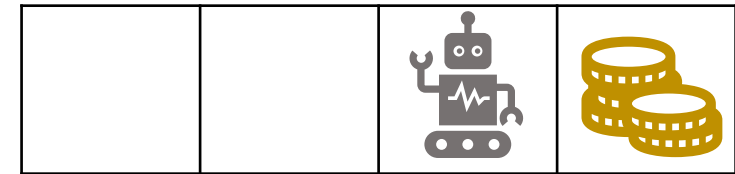
# Tabular Q-Learning

$$Q^{k+1}(s, a) = \mathbb{E}_{s' \sim p(s'|s, a)} \left[ \underbrace{r(s, a, s')}_{= 0} + \gamma \underbrace{\max_{a'} (Q^k(s', a'))}_{= 1} \right]$$

Iteration 1:

State	$Q$				Action		
					←	□	→
					0	0	0
					0	0	1/2
		1/4	1/2	1			
		1/2	1	1			

Environment



$$\gamma = 1/2$$


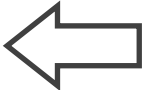

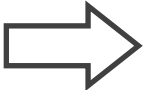


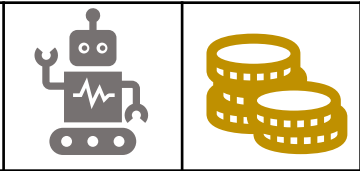
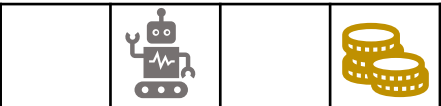
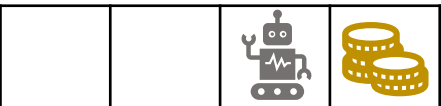
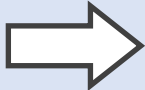

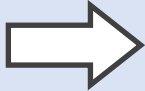

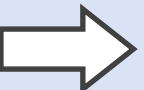




# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$


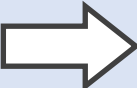
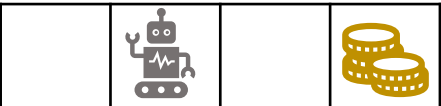
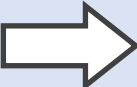
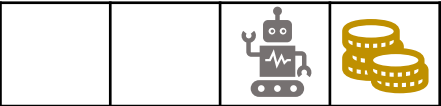
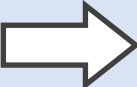


Iteration k:

		Action			Environment			
		$Q$						
State					$\pi$			
		1/8	1/8	1/4	?	$\gamma = 1/2$		
		1/8	1/4	1/2				
		1/4	1/2	1				
		1/2	1	1	 			

# Tabular Q-Learning

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

Iteration k:

		Action			Environment			
		$Q$			$\pi$			
State		1/8	1/8	1/4				
		1/8	1/4	1/2				
		1/4	1/2	1				
		1/2	1	1				
					$\gamma = 1/2$			



# Tabular Q-Learning

---

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

In tabular setting:

- Every iteration leads to a better Q-function + policy
- Converges to optimal Q-function + policy

Limitations:

- Can only be applied to discrete states and actions
- Need to enumerate over all states and actions every iteration

# Large State Spaces

---

Observation:

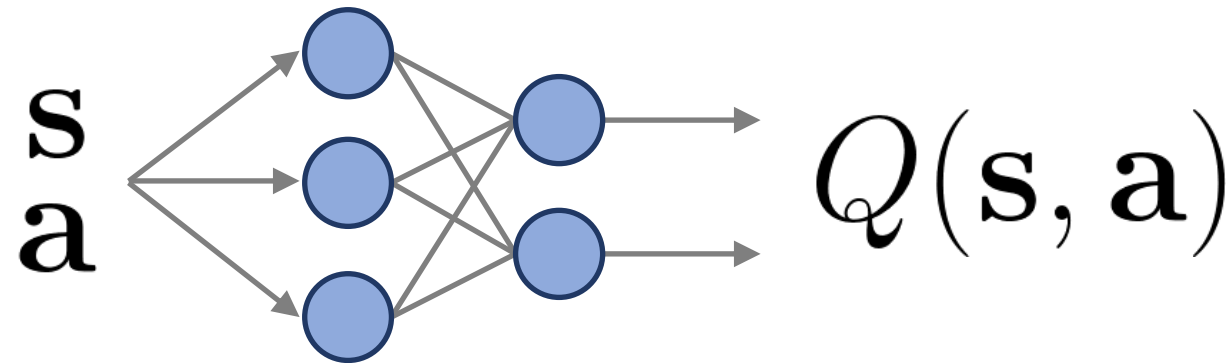
- 64 x 64 image
- 8 bits per pixel
- $2^{8 \times 64 \times 64}$  different states!



# Large/Continuous State Spaces

---

$$\underline{Q^{k+1}(\mathbf{s}, \mathbf{a})} = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

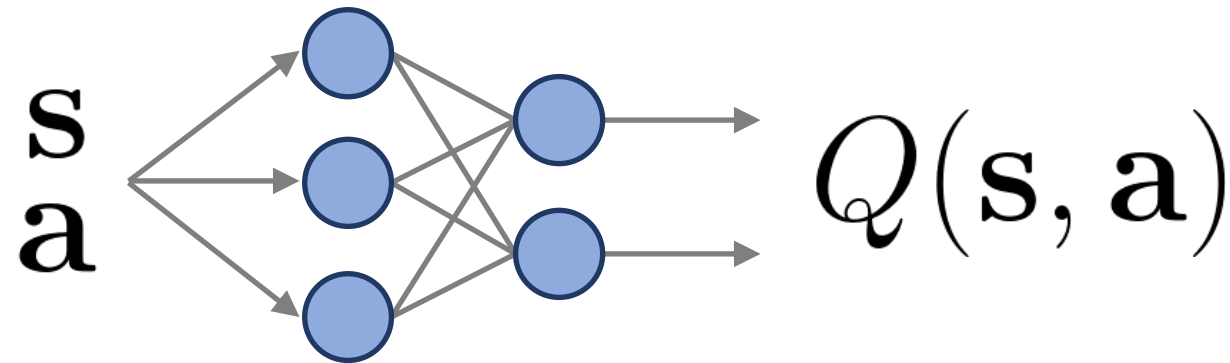


# Large/Continuous State Spaces

---

$$\underline{Q^{k+1}(\mathbf{s}, \mathbf{a})} = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

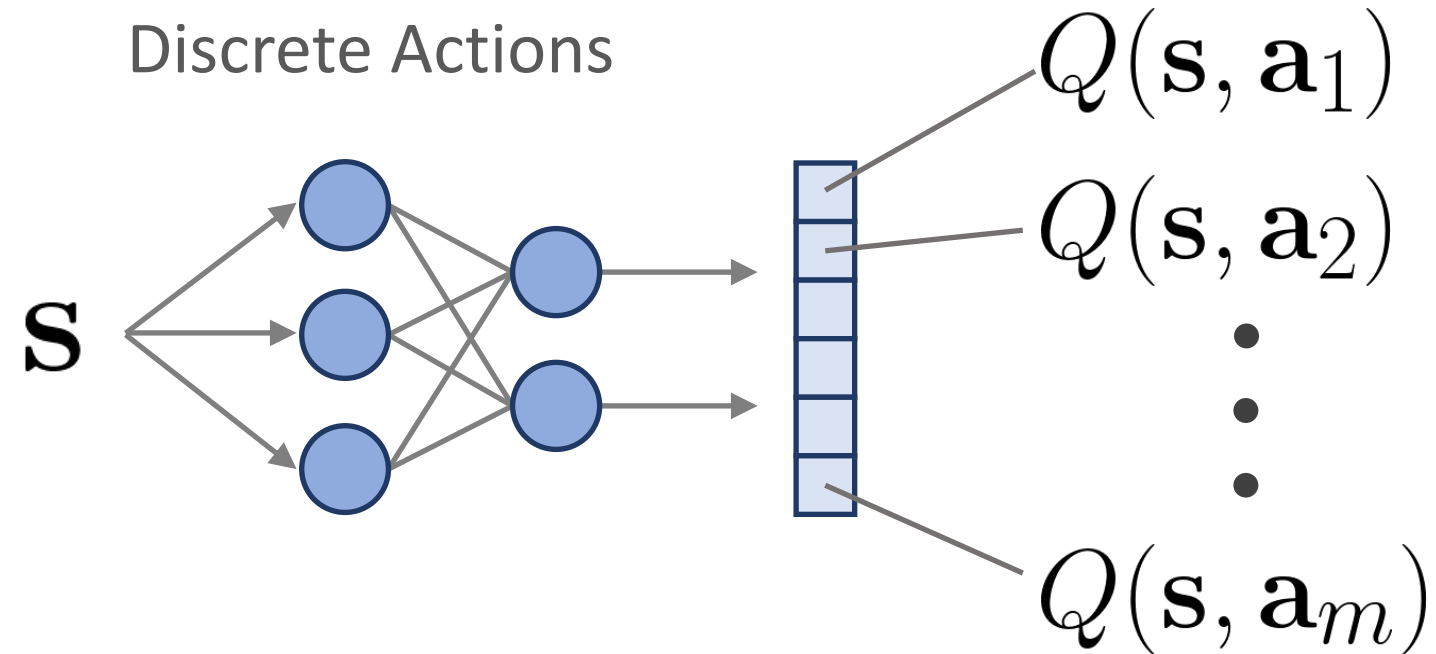
Discrete Actions





# Large/Continuous State Spaces

$$\underline{Q^{k+1}(\mathbf{s}, \mathbf{a})} = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

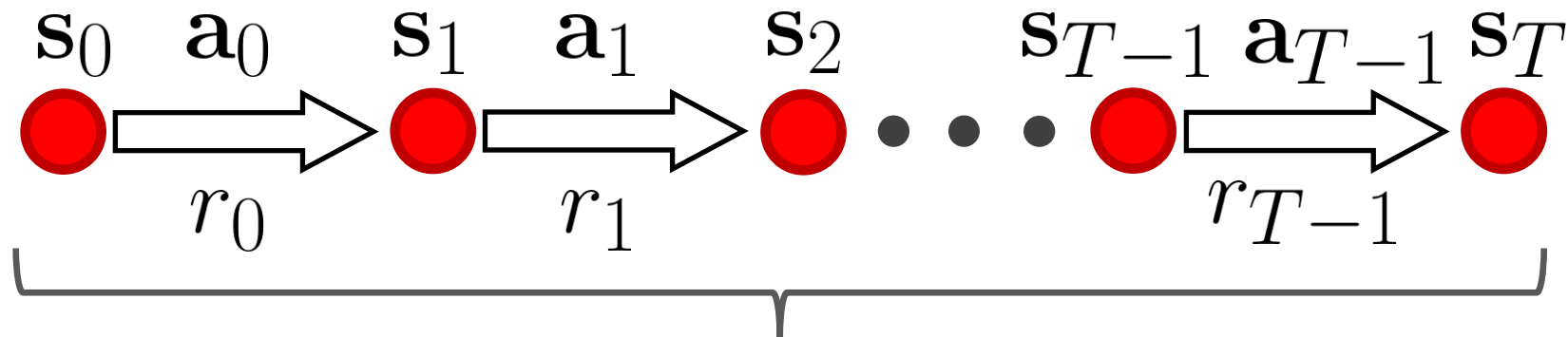


# Large/Continuous State Spaces

---

$$Q^{k+1}(\underline{\mathbf{s}}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

$$Q^k(\mathbf{s}, \mathbf{a}) \Longrightarrow \pi^k(\mathbf{a}|\mathbf{s})$$



$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$$

# Large/Continuous State Spaces

---

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} \left( Q^k(\mathbf{s}', \mathbf{a}') \right) \right]$$

$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$$

Compute target values for each sample  $i$

$$y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$$

Fit new Q-function

$$Q^{k+1} = \arg \min_Q \underbrace{\mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} \left[ (y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2 \right]}_{\text{"Bellman error"}}$$

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  - 9: return  $Q^n$
-



# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  
  - 9: return  $Q^n$
-

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  
  - 9: return  $Q^n$
-

# Q-Learning

---

## ALGORITHM: Q-Learning

---

1:  $Q^0 \leftarrow$  initialize Q-function

2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset

3: **for** iteration  $k = 0, \dots, n - 1$  **do**

4: Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$

How to sample trajectories?



5: Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$

6: Calculate target values for each sample  $i$ :

$$y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$$

7: Update Q-function:

$$Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$$

8: **end for**

9: return  $Q^n$

---

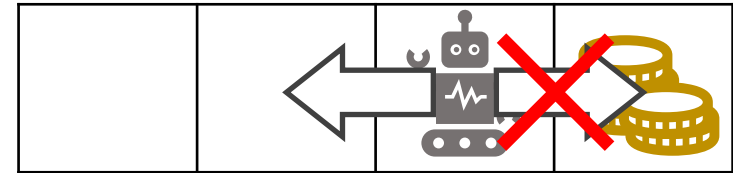
# Sampling

---

$$Q^k(\mathbf{s}, \mathbf{a}) \Longrightarrow \pi^k(\mathbf{a}|\mathbf{s})$$

$$\pi^k(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^k(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$

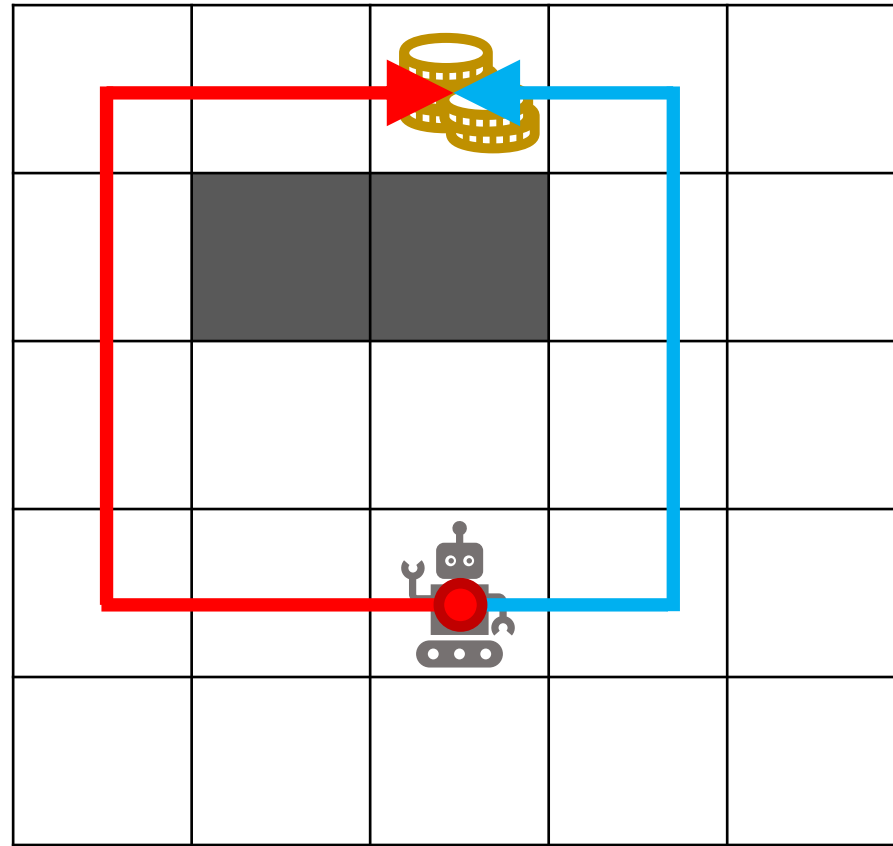
**✗** Insufficient exploration



# Exploration-Exploitation

---

Need to try new actions in case they are better



# Exploration-Exploitation

---

Need to try new actions in case they are better



Keep going to the same restaurant



Try new restaurant

# Exploration-Exploitation

---

Need to try new actions in case they are better

$$\pi^k(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^k(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$



# Exploration-Exploitation

---

Need to try new actions in case they are better

$$\pi^k(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 - \epsilon & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^k(\mathbf{s}, \mathbf{a}') \\ \epsilon & \text{otherwise} \end{cases}$$

Epsilon-greedy exploration

- With probability  $1 - \epsilon$  exploit current best action
- With probability  $\epsilon$  explore new action by sampling a random action
- Start with  $\epsilon = 1$  and then anneal to lower value (e.g.  $\epsilon \rightarrow 0.1$ )

# Boltzmann Exploration

---

Probability of an action is proportion to its “goodness”

$$\pi^k(\mathbf{a}|\mathbf{s}) = \frac{1}{Z} \exp \left( \frac{1}{\beta} Q^k(\mathbf{s}, \mathbf{a}) \right)$$

where,

temperature parameters:  $\beta \in \mathbb{R}$

normalization factor:  $Z = \sum_{\mathbf{a}'} \exp \left( \frac{1}{\beta} Q^k(\mathbf{s}, \mathbf{a}') \right)$

# Testing

---

After training, test with greedy policy

$$\pi^k(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^k(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$

# Q-Learning

---

---

**ALGORITHM: Q-Learning**

---

- 1:  $Q^0 \leftarrow$  initialize Q-function
  - 2:  $\mathcal{D} \leftarrow \{\emptyset\}$  initialize dataset
  
  - 3: **for** iteration  $k = 0, \dots, n - 1$  **do**
  - 4:   Sample trajectory  $\tau$  according to  $Q^k(\mathbf{s}, \mathbf{a})$
  - 5:   Add transitions to dataset  $\mathcal{D} = \mathcal{D} \cup \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$
  
  - 6:   Calculate target values for each sample  $i$ :  
     $y_i = r_i + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}'_i, \mathbf{a}')$
  
  - 7:   Update Q-function:  
     $Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i) \sim \mathcal{D}} [(y_i - Q(\mathbf{s}_i, \mathbf{a}_i))^2]$
  - 8: **end for**
  
  - 9: return  $Q^n$
-

# Q-Learning with Function Approximators

- No improvement guarantees

$$Q^{k+1}(\mathbf{s}, \mathbf{a}) \not\geq Q^k(\mathbf{s}, \mathbf{a}) \qquad J(\pi^{k+1}) \not\geq J(\pi^k)$$

- No convergence guarantees

$$Q^k \not\Rightarrow Q^*$$

- But in practice, it works!

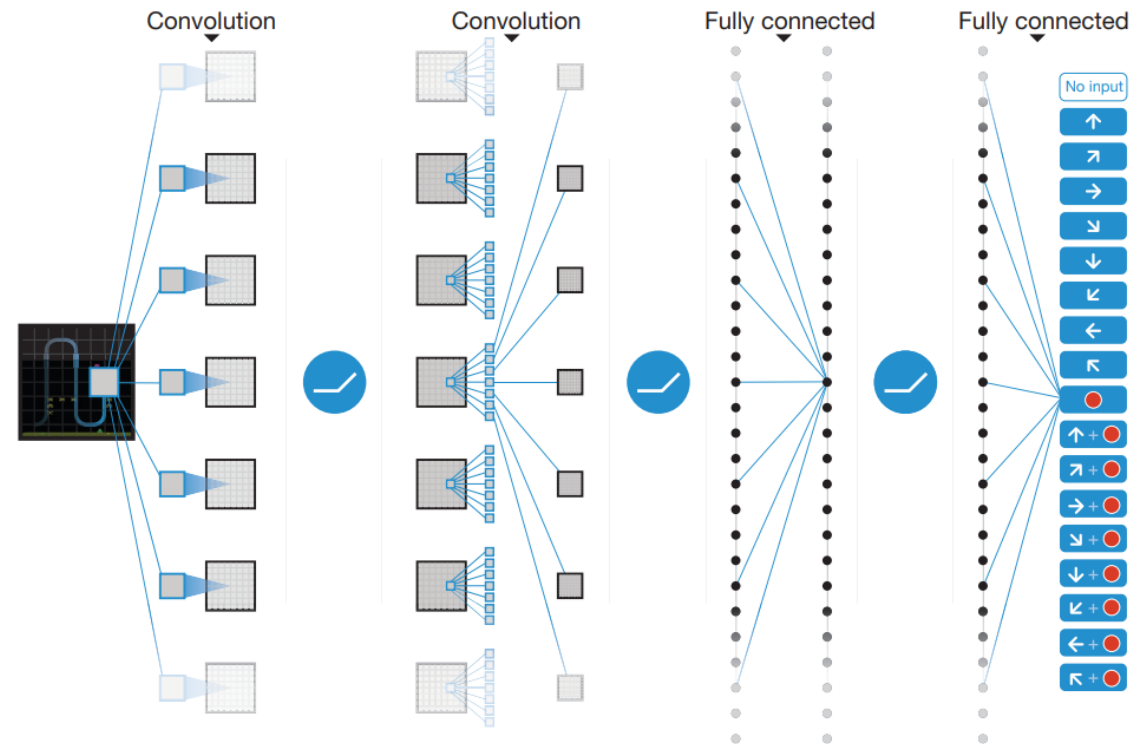
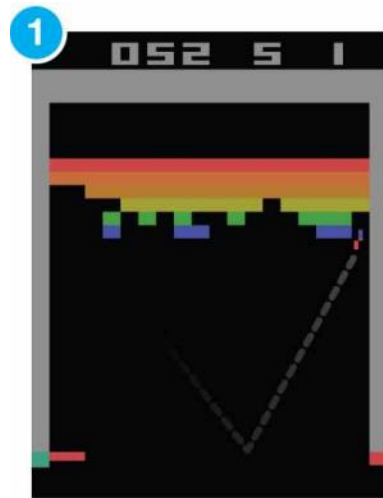
# Deep Q-Networks (DQN)



Human-Level Control Through Deep Reinforcement Learning  
[Mnih et al. 2015]

# Deep Q-Networks (DQN)

Input: 84 x 84 images



Output: Q-values for joystick controls



Human-Level Control Through Deep Reinforcement Learning  
[Mnih et al. 2015]

# Deep Q-Networks (DQN)



Human-Level Control Through Deep Reinforcement Learning  
[Mnih et al. 2015]



# Q-Learning

---

- ✓ Often much more sample efficient than policy gradient
- ✓ Off-policy learning
- ✗ Limited to relatively small discrete action spaces
- ✗ Does not directly optimize performance
  - Lower Bellman error  $\neq$  better performance
- ✗ No convergence guarantees with function approximators

$$Q^{k+1} = \arg \min_Q \mathbb{E}_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{D}} \left[ \left( \left( r + \gamma \max_{\mathbf{a}'} Q^k(\mathbf{s}', \mathbf{a}') \right) - Q(\mathbf{s}, \mathbf{a}) \right)^2 \right]$$

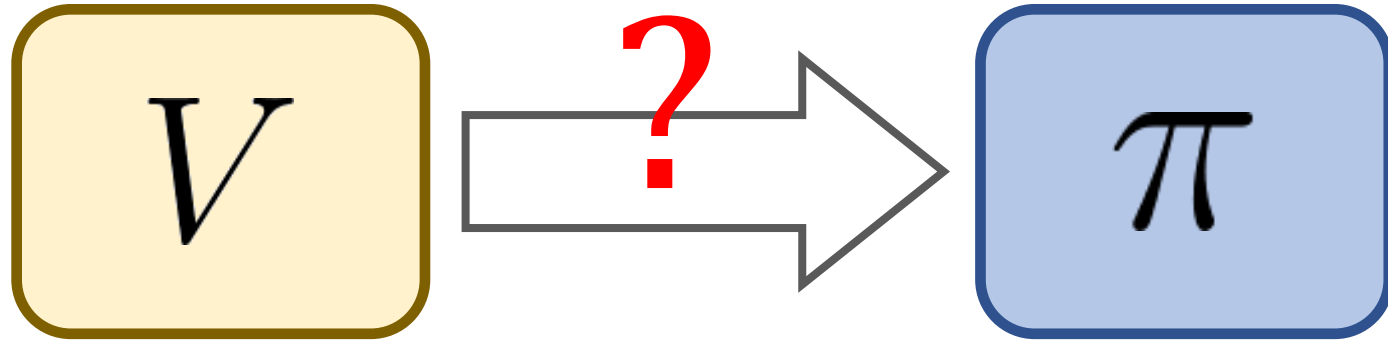
Intractable in large/continuous  
action spaces

# Value Functions

---

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q(\mathbf{s}, \mathbf{a}') \\ 0 & \text{otherwise} \end{cases}$$

What about  $V(\mathbf{s})$ ?



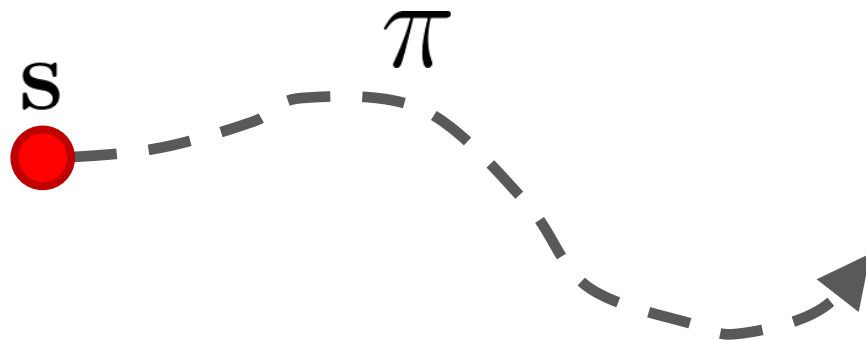
# Value Functions

---

## Value Function

“State Value Function”

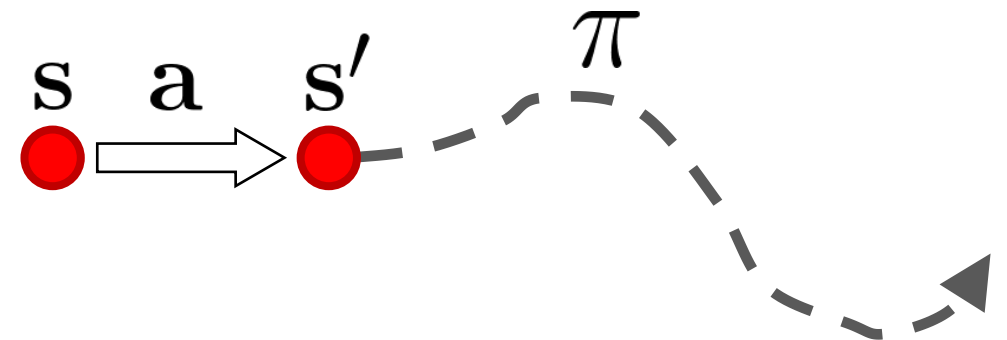
$$V^{\pi}(\mathbf{s}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$



## Q-Function

“State-Action Value Function”

$$Q^{\pi}(\mathbf{s}, \underline{\mathbf{a}}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$



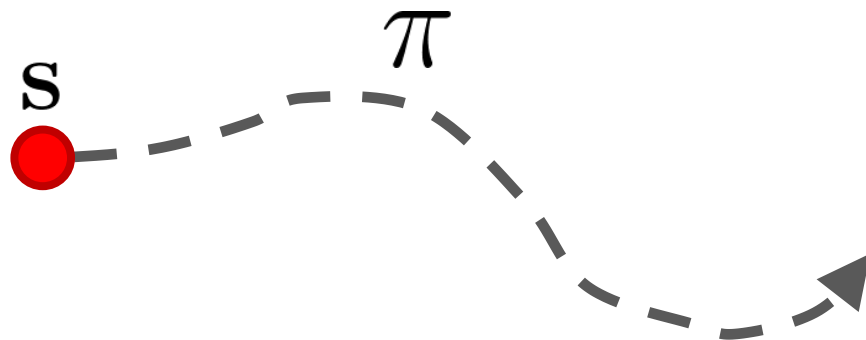
# Value Functions

---

## Value Function

“State Value Function”

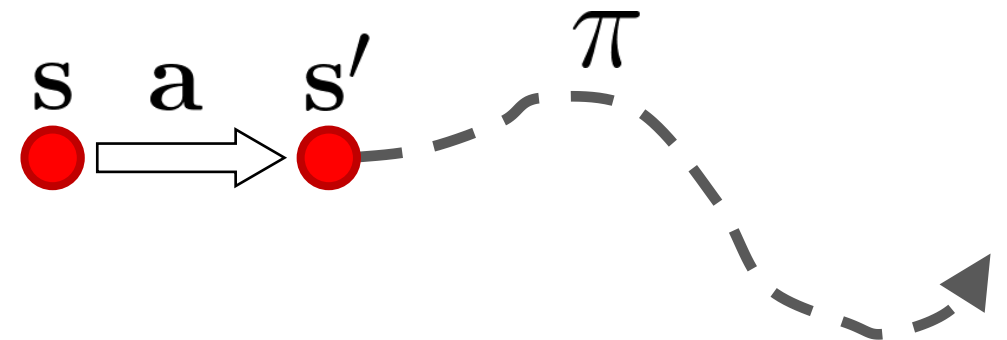
$$V^{\pi}(\underline{\mathbf{s}}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$



## Q-Function

“State-Action Value Function”

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\tau \sim p(\tau | \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a})} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$



# Recursive definition

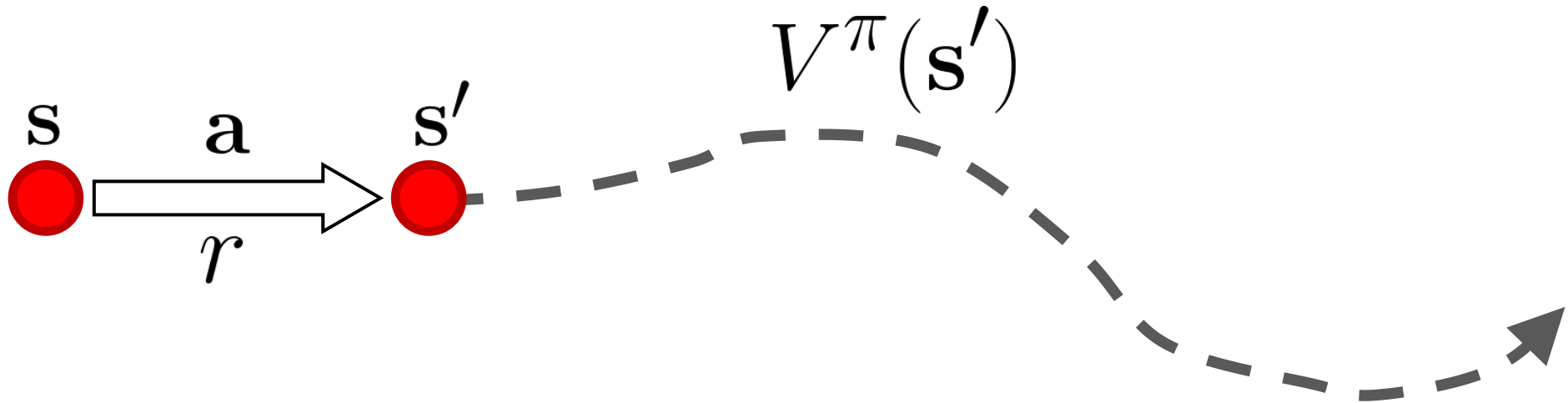
---

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right]$$

# Recursive definition

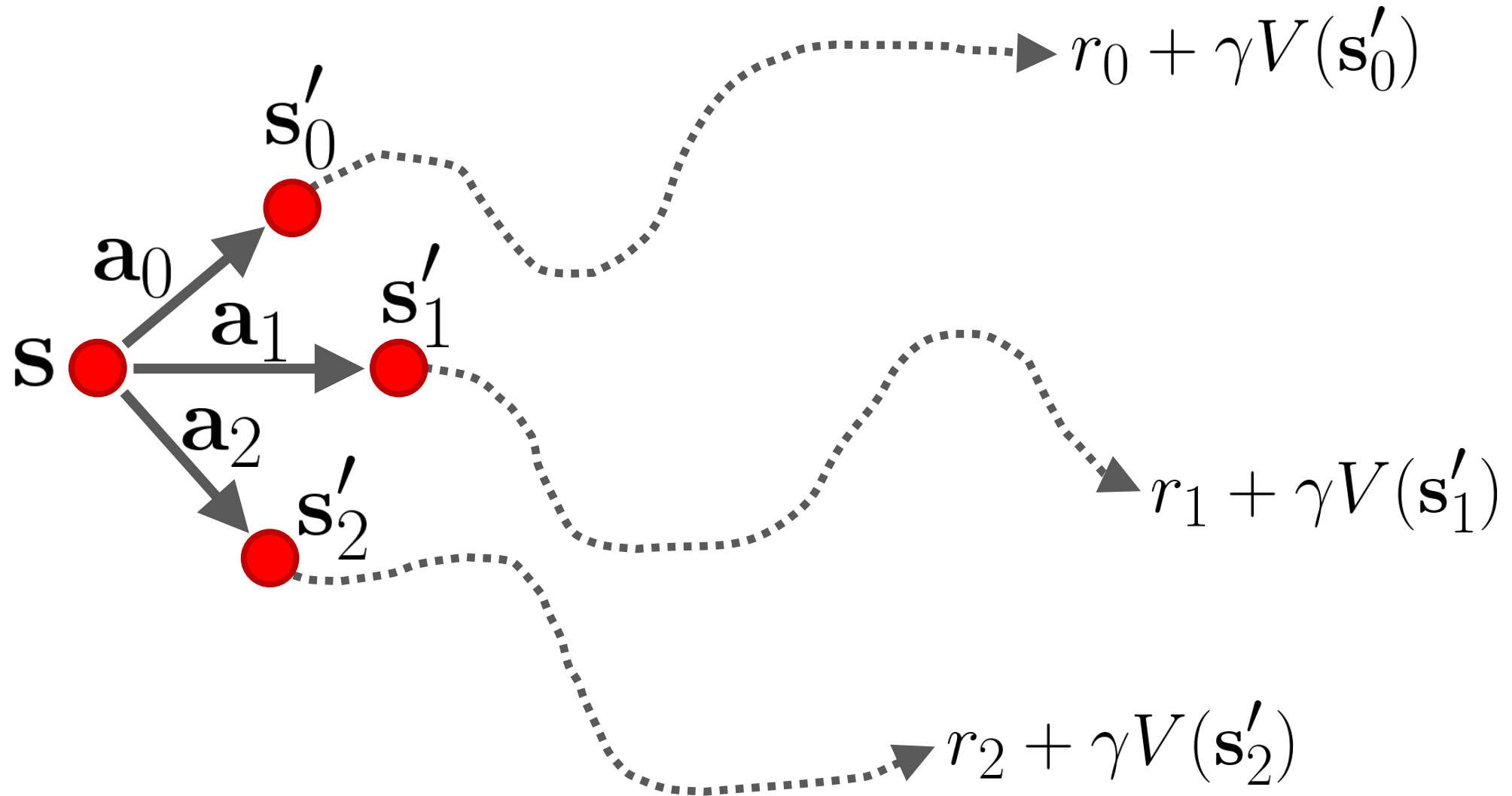
---

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] \right] \\ &= \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} \left[ r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \underline{V^\pi(\mathbf{s}')} \right] \end{aligned}$$



# Value Function

---



# Value Function

---

Value-function:

$$\arg \max_{\mathbf{a}} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V(\mathbf{s}')]$$

Need access to  
dynamics



# Value Function

---

Value-function:

$$\arg \max_{\mathbf{a}} \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V(\mathbf{s}')] ]$$

Q-function:

$$\arg \max_{\mathbf{a}} \underline{Q(\mathbf{s}, \mathbf{a})}$$

Do not need  
dynamics

# Summary

---

- Q-Function
- Q-Learning
- Exploration