

# Reinforcement Learning

Shan-Hung Wu

*shwu@cs.nthu.edu.tw*

Department of Computer Science,  
National Tsing Hua University, Taiwan

Machine Learning

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Outline

## ① Introduction

## ② Markov Decision Process

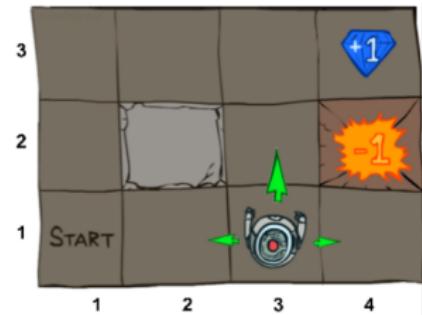
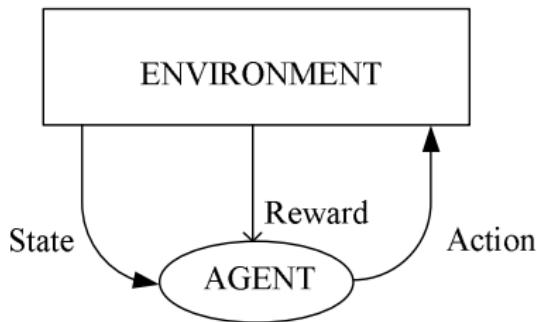
- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

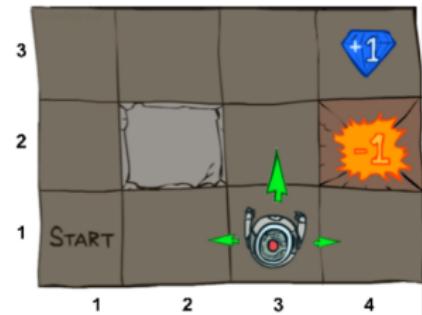
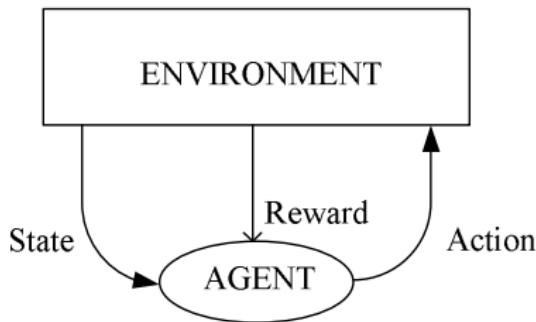
# Reinforcement Learning I

- An agent sees **states**  $s^{(t)}$ 's of an environment, takes **actions**  $a^{(t)}$ 's, and receives **rewards**  $R^{(t)}$ 's (or penalties)



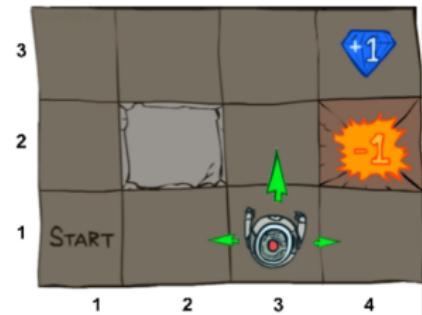
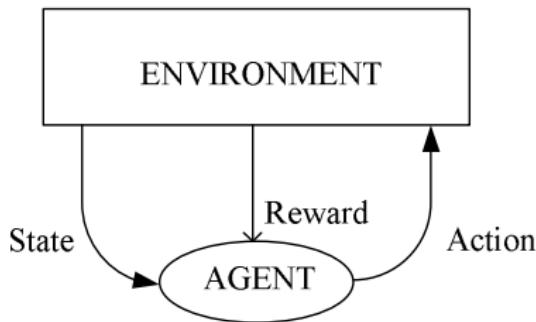
# Reinforcement Learning I

- An agent sees **states**  $s^{(t)}$ 's of an environment, takes **actions**  $a^{(t)}$ 's, and receives **rewards**  $R^{(t)}$ 's (or penalties)
  - Environment does **not** change over time



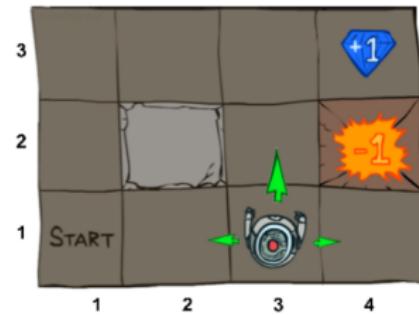
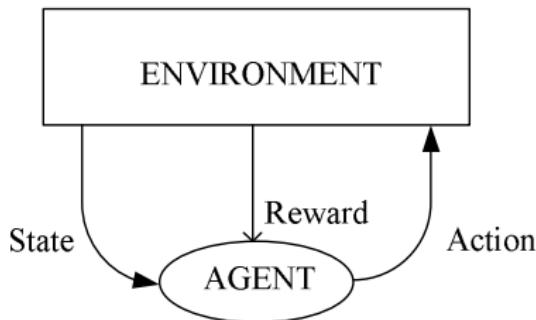
# Reinforcement Learning I

- An agent sees **states**  $s^{(t)}$ 's of an environment, takes **actions**  $a^{(t)}$ 's, and receives **rewards**  $R^{(t)}$ 's (or penalties)
  - Environment does **not** change over time
  - The state of the environment may change due to an action



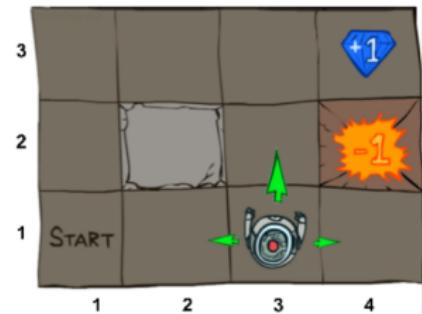
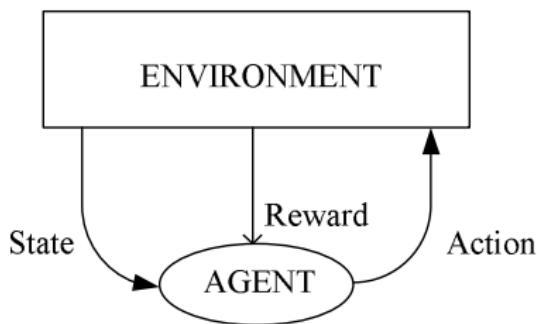
# Reinforcement Learning I

- An agent sees **states**  $s^{(t)}$ 's of an environment, takes **actions**  $a^{(t)}$ 's, and receives **rewards**  $R^{(t)}$ 's (or penalties)
  - Environment does **not** change over time
  - The state of the environment may change due to an action
  - Reward  $R^{(t)}$  may depend on  $s^{(t+1)}, s^{(t)}, \dots$  or  $a^{(t)}, a^{(t-1)}, \dots$



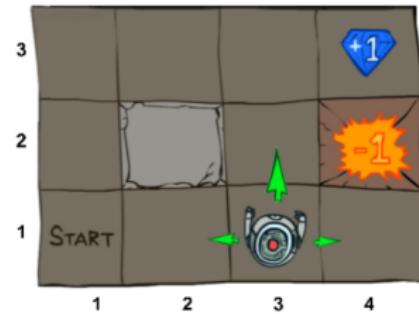
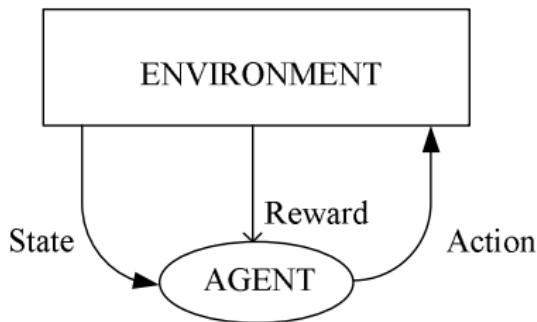
# Reinforcement Learning II

- Goal: to learn the best **policy**  $\pi^*(s^{(t)}) = a^{(t)}$  that maximizes the **total** reward  $\sum_t R^{(t)}$



# Reinforcement Learning II

- Goal: to learn the best **policy**  $\pi^*(s^{(t)}) = a^{(t)}$  that maximizes the **total** reward  $\sum_t R^{(t)}$
- Training:
  - 1 Perform trial-and-error runs
  - 2 Learn from the experience



# Compared to Supervised/Unsupervised Learning

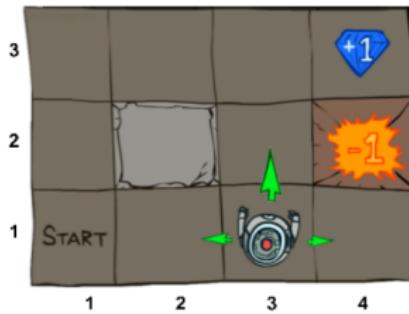
- $\pi^*(s^{(t)}) = \mathbf{a}^{(t)}$  maximizing total reward  $\sum_t R^{(t)}$  vs.  $f^*(\mathbf{x}^{(i)}) = \mathbf{y}^{(i)}$   
minimizing total loss

# Compared to Supervised/Unsupervised Learning

- $\pi^*(s^{(t)}) = \mathbf{a}^{(t)}$  maximizing total reward  $\sum_t R^{(t)}$  vs.  $f^*(\mathbf{x}^{(i)}) = \mathbf{y}^{(i)}$  minimizing total loss
- Examples  $\mathbf{x}^{(i)}$ 's are i.i.d., but not in RL
  - $s^{(t+1)}$  may depend on  $s^{(t)}, s^{(t-1)}, \dots$  and  $\mathbf{a}^{(t)}, \mathbf{a}^{(t-1)}, \dots$

# Compared to Supervised/Unsupervised Learning

- $\pi^*(s^{(t)}) = a^{(t)}$  maximizing total reward  $\sum_t R^{(t)}$  vs.  $f^*(x^{(i)}) = y^{(i)}$  minimizing total loss
- Examples  $x^{(i)}$ 's are i.i.d., but not in RL
  - $s^{(t+1)}$  may depend on  $s^{(t)}, s^{(t-1)}, \dots$  and  $a^{(t)}, a^{(t-1)}, \dots$
- No **what** to predict ( $y^{(i)}$ 's), just **how good** a prediction is ( $R^{(t)}$ 's)
  - $R^{(t)}$ 's are also called the **critics**



# Applications

- Sequential decision making and control problems



Kohl and Stone, 2004



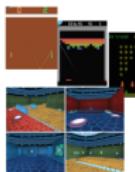
Ng et al, 2004



Tedrake et al, 2005

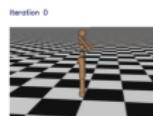


Kober and Peters, 2009



Mnih et al 2013 (DQN)  
Mnih et al, 2015 (A3C)

Silver et al, 2014 (DPG)  
Lillicrap et al, 2015 (DDPG)



Schulman et al,  
2016 (TRPO + GAE)



Levine\*, Finn\*, et  
al, 2016  
(GPS)



Silver\*, Huang\*, et  
al, 2016  
(AlphaGo)

# Applications

- Sequential decision making and control problems



Kohl and Stone, 2004



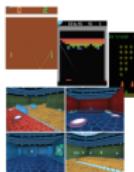
Ng et al, 2004



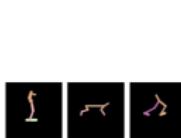
Tedrake et al, 2005



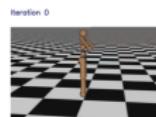
Kober and Peters, 2009



Mnih et al 2013 (DQN)  
Mnih et al, 2015 (A3C)



Silver et al, 2014 (DPG)  
Lillicrap et al, 2015 (DDPG)



Schulman et al,  
2016 (TRPO + GAE)



Levine\*, Finn\*, et  
al, 2016  
(GPS)



Silver\*, Huang\*, et  
al, 2016  
(AlphaGo)

- From machine learning to AI that changes the world



# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Markov Processes

- A **random process**  $\{s^{(t)}\}_t$  is a collection of time-indexed random variables
  - A classic way to model dependency between input samples  $\{s^{(t)}\}_t$
- A random process is called a **Markov process** if it satisfies the **Markov property**:

$$P(s^{(t+1)} | s^{(t)}, s^{(t-1)}, \dots) = P(s^{(t+1)} | s^{(t)})$$

# Markov Processes

- A **random process**  $\{s^{(t)}\}_t$  is a collection of time-indexed random variables
  - A classic way to model dependency between input samples  $\{s^{(t)}\}_t$
- A random process is called a **Markov process** if it satisfies the **Markov property**:

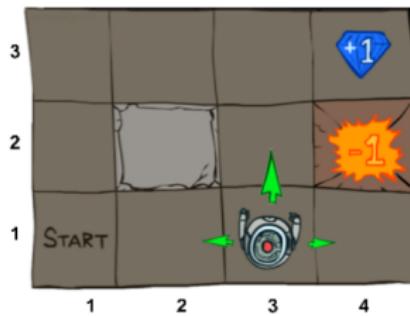
$$P(s^{(t+1)} | s^{(t)}, s^{(t-1)}, \dots) = P(s^{(t+1)} | s^{(t)})$$

- For those who knows Markov process:\*\*

	States are fully observable	States are partially observable
Transition is autonomous	Markov chains	Hidden Markov models
Transition is controlled	<b>Markov decision processes (MDP)</b>	Partially observable MDP

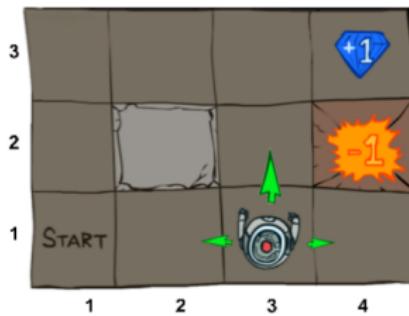
# Markov Decision Process I

- A Markov decision process (MDP) is defined by
  - $\mathbb{S}$  the state space;  $\mathbb{A}$  the action space
  - Start state  $s^{(0)}$
  - $P(s'|s; a)$  the **transition distribution** controlled by actions; fixed over time  $t$
  - $R(s, a, s') \in \mathbb{R}$  (or simply  $R(s')$ ) the deterministic reward function
  - $\gamma \in [0, 1]$  is the **discount factor**
  - $H \in \mathbb{N}$  the **horizon**; can be infinite



# Markov Decision Process I

- A Markov decision process (MDP) is defined by
  - $\mathbb{S}$  the state space;  $\mathbb{A}$  the action space
  - Start state  $s^{(0)}$
  - $P(s'|s; a)$  the **transition distribution** controlled by actions; fixed over time  $t$
  - $R(s, a, s') \in \mathbb{R}$  (or simply  $R(s')$ ) the deterministic reward function
  - $\gamma \in [0, 1]$  is the **discount factor**
  - $H \in \mathbb{N}$  the **horizon**; can be infinite
- An absorbing/terminal state transit to itself with probability 1



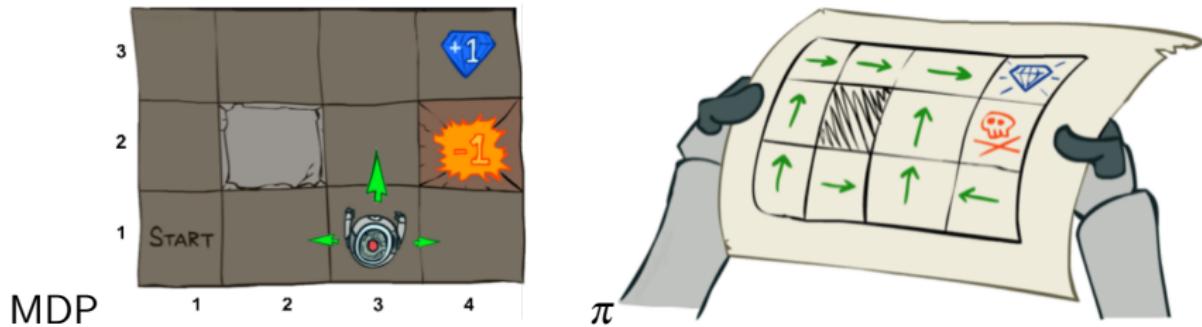
# Markov Decision Process II

- Given a policy  $\pi(s) = a$ , an MDP proceeds as follows:

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)},$$

with the accumulative reward

$$R(s^{(0)}, a^{(0)}, s^{(1)}) + \gamma R(s^{(1)}, a^{(1)}, s^{(2)}) + \dots + \gamma^{H-1} R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$



# Markov Decision Process II

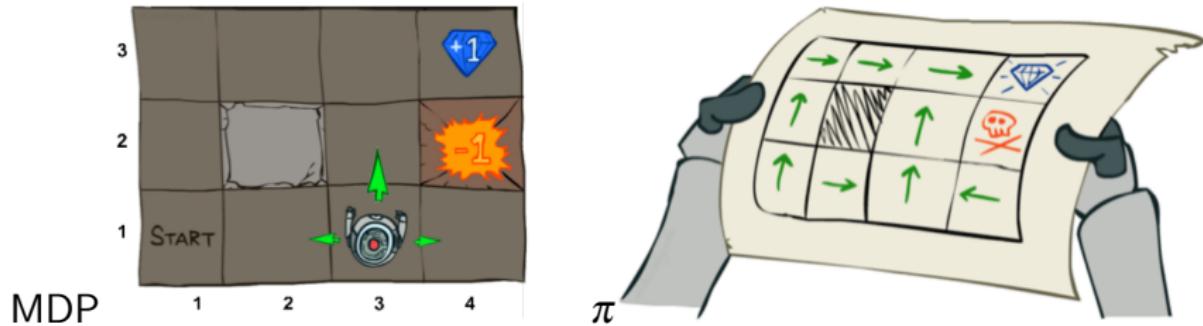
- Given a policy  $\pi(s) = a$ , an MDP proceeds as follows:

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)},$$

with the accumulative reward

$$R(s^{(0)}, a^{(0)}, s^{(1)}) + \gamma R(s^{(1)}, a^{(1)}, s^{(2)}) + \dots + \gamma^{H-1} R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$

- To accrue rewards as soon as possible (prefer a short path)*
- Different accumulative rewards in different trials*



# Goal

- Given a policy  $\pi$ , the expected accumulative reward collected by taking actions following  $\pi$  can be express by:

$$V_\pi = \mathbb{E}_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left( \sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- Goal: to find the optimal policy

$$\pi^* = \arg \max_{\pi} V_{\pi}$$

# Goal

- Given a policy  $\pi$ , the expected accumulative reward collected by taking actions following  $\pi$  can be express by:

$$V_\pi = \mathbb{E}_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left( \sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- Goal: to find the optimal policy

$$\pi^* = \arg \max_{\pi} V_{\pi}$$

- How?

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Optimal Value Function

$$\pi^* = \arg \max_{\pi} E_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left( \sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

# Optimal Value Function

$$\pi^* = \arg \max_{\pi} E_{s^{(0)}, \dots, s^{(H)}} \left( \sum_{t=0}^H \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}); \pi \right)$$

- *Optimal value function:*

$$V^{*(h)}(s) = \max_{\pi} E_{s^{(1)}, \dots, s^{(h)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- Maximum expected accumulative reward when starting from state  $s$  and acting optimally for  $h$  steps

# Optimal Value Function

$$\pi^* = \arg \max_{\pi} E_{s^{(0)}, \dots, s^{(H)}} \left( \sum_{t=0}^H \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}); \pi \right)$$

- *Optimal value function:*

$$V^{*(h)}(s) = \max_{\pi} E_{s^{(1)}, \dots, s^{(h)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- Maximum expected accumulative reward when starting from state  $s$  and acting optimally for  $h$  steps
- Having  $V^{*(H-1)}(s)$  for each  $s$ , we can solve  $\pi^*$  easily by

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^{*(H-1)}(s')], \forall s$$

in  $O(|\mathbb{S}||\mathbb{A}|)$  time

# Optimal Value Function

$$\pi^* = \arg \max_{\pi} E_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left( \sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- *Optimal value function:*

$$V^{*(h)}(\mathbf{s}) = \max_{\pi} E_{\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(h)}} \left( \sum_{t=0}^h \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) | \mathbf{s}^{(0)} = s; \pi \right)$$

- Maximum expected accumulative reward when starting from state  $s$  and acting optimally for  $h$  steps
- Having  $V^{*(H-1)}(s)$  for each  $s$ , we can solve  $\pi^*$  easily by

$$\pi^*(s) = \arg \max_{\mathbf{a}} \sum_{s'} P(s'|s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-1)}(s')], \forall s$$

in  $O(|\mathbb{S}||\mathbb{A}|)$  time

- How to obtain  $V^{*(H-1)}(s)$  for each  $s$ ?

# Dynamic Programming

$$V^{*(h)}(s) = \max_{\pi} E_{s^{(1)}, \dots, s^{(H)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

# Dynamic Programming

$$V^{*(h)}(s) = \max_{\pi} \mathbb{E}_{s^{(1)}, \dots, s^{(H)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- $h = H - 1$ :

$$V^{*(H-1)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-2)}(s')], \forall s$$

# Dynamic Programming

$$V^{*(h)}(s) = \max_{\pi} \mathbb{E}_{s^{(1)}, \dots, s^{(H)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- $h = H - 1$ :

$$V^{*(H-1)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-2)}(s')], \forall s$$

- $h = H - 2$ :

$$V^{*(H-2)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-3)}(s')], \forall s$$

# Dynamic Programming

$$V^{*(h)}(s) = \max_{\pi} \mathbb{E}_{s^{(1)}, \dots, s^{(H)}} \left( \sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- $h = H - 1$ :

$$V^{*(H-1)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-2)}(s')], \forall s$$

- $h = H - 2$ :

$$V^{*(H-2)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(H-3)}(s')], \forall s$$

- $h = 0$ :

$$V^{*(0)}(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^{*(-1)}(s')], \forall s$$

- $h = -1$ :

$$V^{*(-1)}(s) = 0, \forall s$$

# Algorithm: Value Iteration (Finite Horizon)

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H)$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $V^*(s) \leftarrow 0$ ;

**for**  $h \leftarrow 0$  **to**  $H - 1$  **do**

**foreach**  $s$  **do**

$| V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')];$

**end**

**end**

**foreach**  $s$  **do**

$| \pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')];$

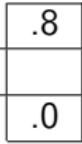
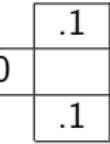
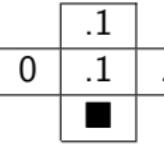
**end**

**Algorithm 1:** Value iteration with finite horizon.

# Example

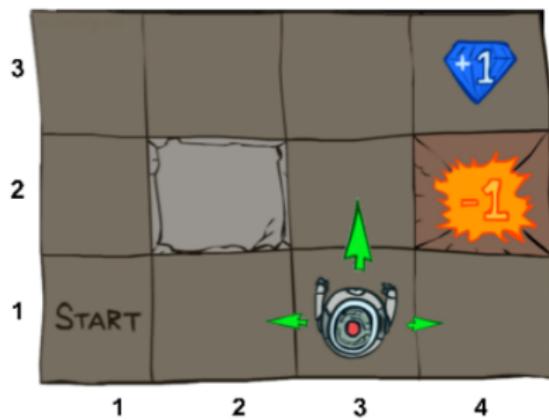
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; \mathbf{a})$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9



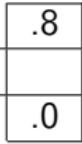
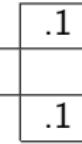
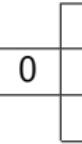
0.00	0.00	0.00	<input type="text"/>
0.00		0.00	<input type="text"/>
0.00	0.00	0.00	0.00

VALUES AFTER 0 ITERATIONS

# Example

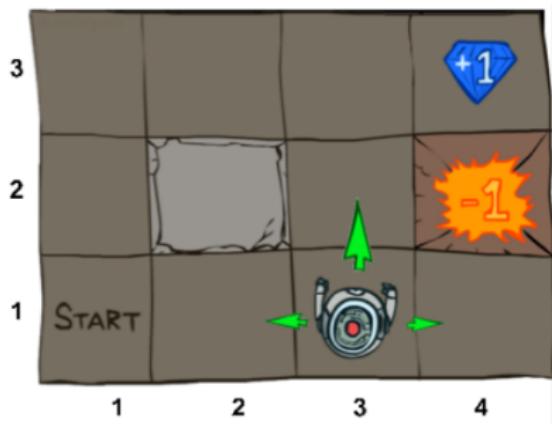
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; \mathbf{a})$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9



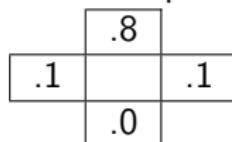
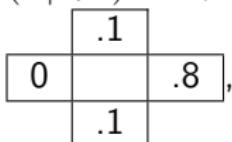
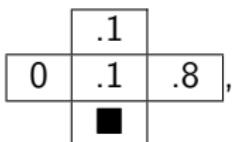
0.00	0.00	0.00	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 1 ITERATIONS

# Example

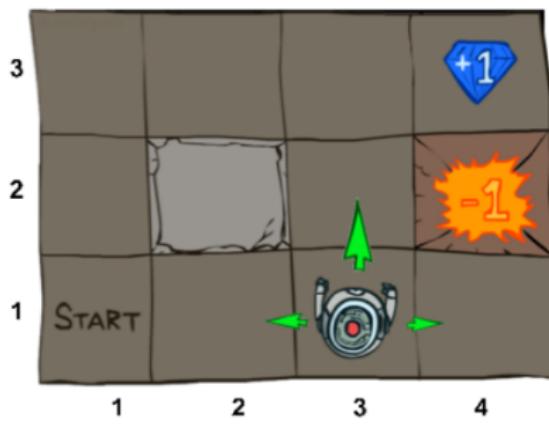
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; \mathbf{a})$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9



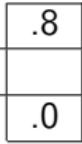
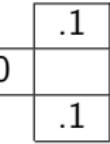
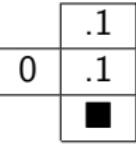
0.00	0.00	0.72	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 2 ITERATIONS

# Example

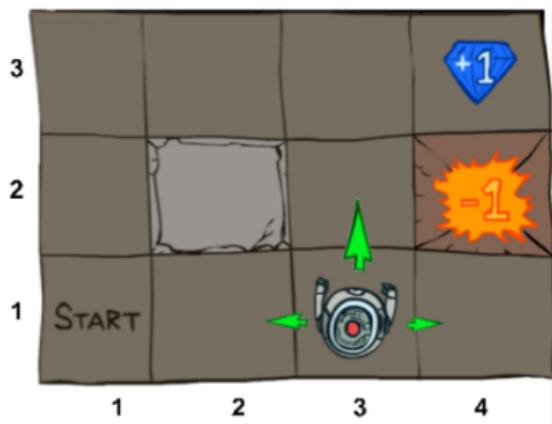
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; a)$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9

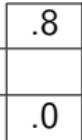
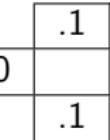
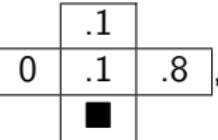


0.00	0.52	0.78	1.00
0.00		0.43	-1.00
0.00	0.00	0.00	0.00

# Example

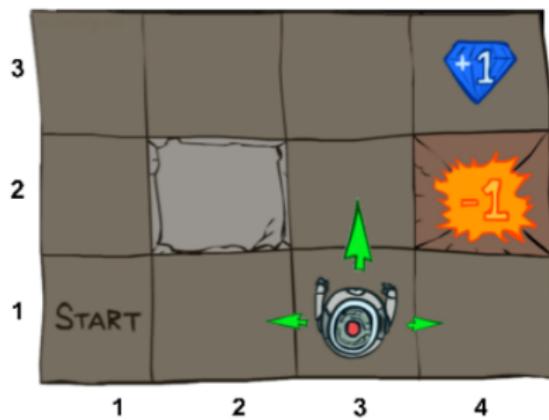
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; a)$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9

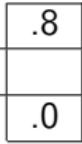
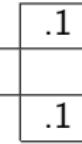
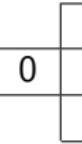


0.37	0.66	0.83	1.00
0.00		0.51	-1.00
0.00	0.00	0.31	0.00

# Example

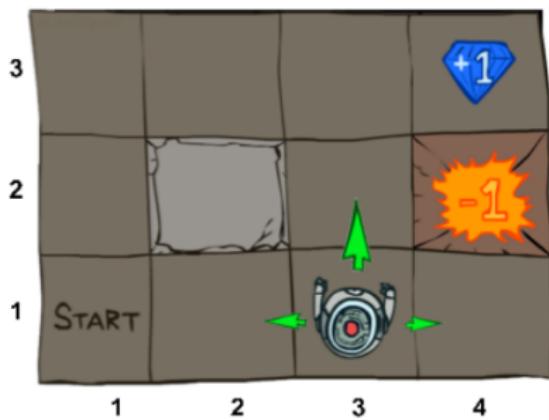
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; \mathbf{a})$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9

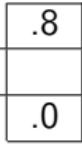
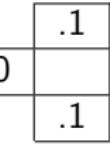
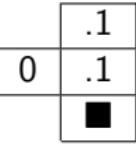


0.51	0.72	0.84	1.00
0.27		0.55	-1.00
0.00	0.22	0.37	0.13

# Example

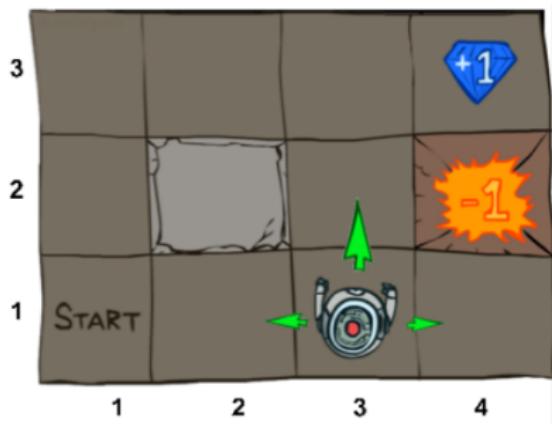
- MDP settings:

- Actions: "up," "down," "left," "right"
- Noise of transition probability  $P(s'|s; a)$ : 0.2,

e.g., up:  right: , right: ,

etc.

- $\gamma$ : 0.9



0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.43	0.48	0.28

VALUES AFTER 100 ITERATIONS

# Infinite Horizon & Bellman Optimality Equation

- Recurrence of optimal values:

$$V^{*(h)}(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^{*(h-1)}(s')], \forall s$$

# Infinite Horizon & Bellman Optimality Equation

- Recurrence of optimal values:

$$V^{*(h)}(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^{*(h-1)}(s')], \forall s$$

- When  $h \rightarrow \infty$ , we have the *Bellman optimality equation*:

$$V^*(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')], \forall s$$

# Infinite Horizon & Bellman Optimality Equation

- Recurrence of optimal values:

$$V^{*(h)}(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^{*(h-1)}(s')], \forall s$$

- When  $h \rightarrow \infty$ , we have the **Bellman optimality equation**:

$$V^*(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')], \forall s$$

- Optimal policy:

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')], \forall s$$

# Infinite Horizon & Bellman Optimality Equation

- Recurrence of optimal values:

$$V^{*(h)}(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^{*(h-1)}(s')], \forall s$$

- When  $h \rightarrow \infty$ , we have the **Bellman optimality equation**:

$$V^*(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')], \forall s$$

- Optimal policy:

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')], \forall s$$

- When  $h \rightarrow \infty$ ,  $\pi^*$  is

- **Stationary**: the optimal action at a state  $s$  is the same at all times (efficient to store)
- **Memoryless**: independent with  $s^{(0)}$

# Algorithm: Value Iteration (Infinite Horizon)

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $V^*(s) \leftarrow 0$ ;

**repeat**

**foreach**  $s$  **do**

$| V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')];$   
    **end**

**until**  $V^*(s)$  's converge;

**foreach**  $s$  **do**

$| \pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')];$

**end**

**Algorithm 2:** Value iteration with infinite horizon.

# Convergence

## Theorem

*Value iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

# Convergence

## Theorem

*Value iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

- Intuition:

$$\begin{aligned} V^*(\mathbf{s}) - V^{*(H)}(\mathbf{s}) &= \gamma^{H+1} R(\mathbf{s}^{(H+1)}, \mathbf{a}^{(H+1)}, \mathbf{s}^{(H+2)}) \\ &\quad + \gamma^{H+2} R(\mathbf{s}^{(H+2)}, \mathbf{a}^{(H+2)}, \mathbf{s}^{(H+3)}) + \dots \\ &\leq \gamma^{H+1} R_{\max} + \gamma^{H+2} R_{\max} + \dots \\ &= \frac{\gamma^{H+1}}{1-\gamma} R_{\max} \end{aligned}$$

- Goes to 0 as  $H \rightarrow \infty$
- Hence,  $V^{*(H)}(\mathbf{s}) \xrightarrow{H \rightarrow \infty} V^*(\mathbf{s})$

# Convergence

## Theorem

*Value iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

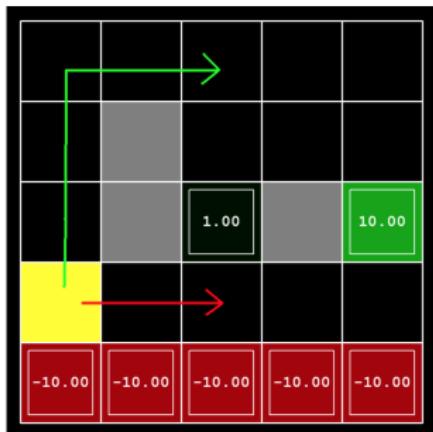
- Intuition:

$$\begin{aligned} V^*(\mathbf{s}) - V^{*(H)}(\mathbf{s}) &= \gamma^{H+1} R(\mathbf{s}^{(H+1)}, \mathbf{a}^{(H+1)}, \mathbf{s}^{(H+2)}) \\ &\quad + \gamma^{H+2} R(\mathbf{s}^{(H+2)}, \mathbf{a}^{(H+2)}, \mathbf{s}^{(H+3)}) + \dots \\ &\leq \gamma^{H+1} R_{\max} + \gamma^{H+2} R_{\max} + \dots \\ &= \frac{\gamma^{H+1}}{1-\gamma} R_{\max} \end{aligned}$$

- Goes to 0 as  $H \rightarrow \infty$
- Hence,  $V^{*(H)}(\mathbf{s}) \xrightarrow{H \rightarrow \infty} V^*(\mathbf{s})$
- Assumed that  $R(\cdot) \geq 0$ ; still holds if rewards can be negative
  - by using  $\max |R(\cdot)|$  and bounding from both sides

# Effect of MDP Parameters

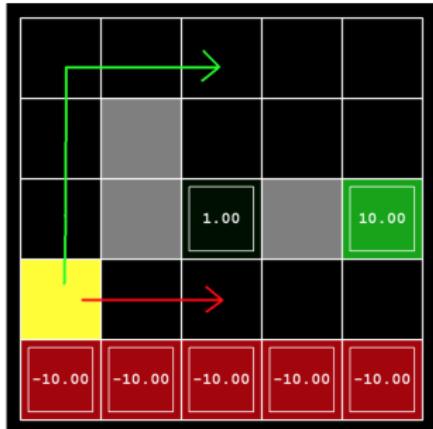
- How does the noise of  $P(s'|s; \mathbf{a})$  and  $\gamma$  affect  $\pi^*$ ?



- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

# Effect of MDP Parameters

- How does the noise of  $P(s'|s; a)$  and  $\gamma$  affect  $\pi^*$ ?

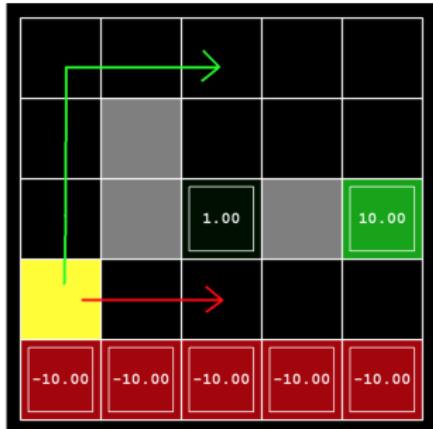


- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

- $\pi^*$  prefers the close exit (+1); risking the cliff (-10)?
- $\pi^*$  prefers the close exit (+1); avoiding the cliff (-10)?
- $\pi^*$  prefers the distant exit (+10); risking the cliff (-10)?
- $\pi^*$  prefers the distant exit (+10); avoiding the cliff (-10)?

# Effect of MDP Parameters

- How does the noise of  $P(s'|s; a)$  and  $\gamma$  affect  $\pi^*$ ?

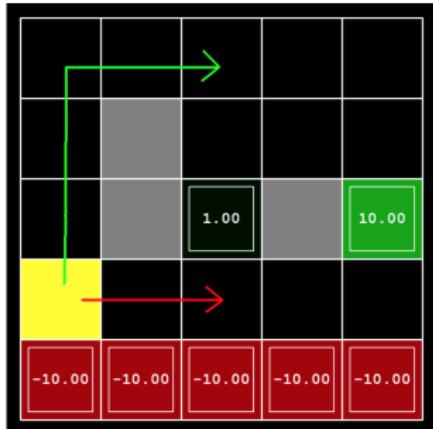


- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

- $\pi^*$  prefers the close exit (+1); risking the cliff (-10)? (4)
- $\pi^*$  prefers the close exit (+1); avoiding the cliff (-10)?
- $\pi^*$  prefers the distant exit (+10); risking the cliff (-10)?
- $\pi^*$  prefers the distant exit (+10); avoiding the cliff (-10)?

# Effect of MDP Parameters

- How does the noise of  $P(s'|s; a)$  and  $\gamma$  affect  $\pi^*$ ?

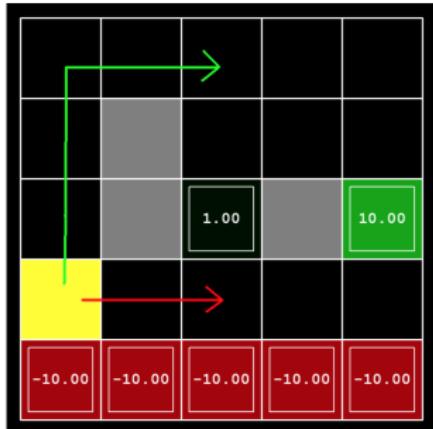


- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

- $\pi^*$  prefers the close exit (+1); risking the cliff (-10)? (4)
- $\pi^*$  prefers the close exit (+1); avoiding the cliff (-10)? (3)
- $\pi^*$  prefers the distant exit (+10); risking the cliff (-10)?
- $\pi^*$  prefers the distant exit (+10); avoiding the cliff (-10)?

# Effect of MDP Parameters

- How does the noise of  $P(s'|s; a)$  and  $\gamma$  affect  $\pi^*$ ?

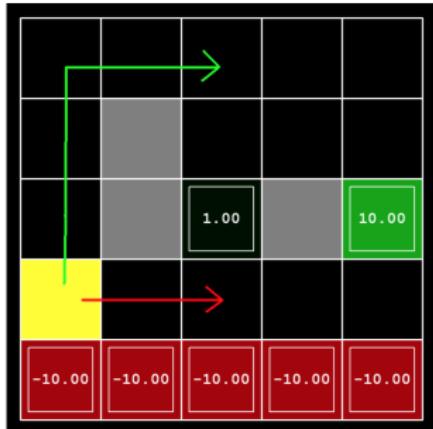


- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

- $\pi^*$  prefers the close exit (+1); risking the cliff (-10)? (4)
- $\pi^*$  prefers the close exit (+1); avoiding the cliff (-10)? (3)
- $\pi^*$  prefers the distant exit (+10); risking the cliff (-10)? (2)
- $\pi^*$  prefers the distant exit (+10); avoiding the cliff (-10)?

# Effect of MDP Parameters

- How does the noise of  $P(s'|s; a)$  and  $\gamma$  affect  $\pi^*$ ?



- ①  $\gamma = 0.99$ , noise = 0.5
- ②  $\gamma = 0.99$ , noise = 0
- ③  $\gamma = 0.1$ , noise = 0.5
- ④  $\gamma = 0.1$ , noise = 0

- $\pi^*$  prefers the close exit (+1); risking the cliff (-10)? (4)
- $\pi^*$  prefers the close exit (+1); avoiding the cliff (-10)? (3)
- $\pi^*$  prefers the distant exit (+10); risking the cliff (-10)? (2)
- $\pi^*$  prefers the distant exit (+10); avoiding the cliff (-10)? (1)

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Goal

- Given an MDP  $(\mathbb{S}, \mathbb{A}, \mathbf{P}, R, \gamma, H \rightarrow \infty)$
- Expected accumulative reward collected by taking actions following a policy  $\pi$ :

$$V_\pi = \mathbb{E}_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left( \sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- Goal: to find the optimal policy

$$\pi^* = \arg \max_{\pi} V_{\pi}$$

# Algorithm: Policy Iteration (Simplified)

- Given a  $\pi$ , define its **value function**:

$$V_\pi(s) = \mathbb{E}_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- Expected accumulative reward when starting from state  $s$  and acting based on  $\pi$

# Algorithm: Policy Iteration (Simplified)

- Given a  $\pi$ , define its **value function**:

$$V_\pi(s) = \mathbb{E}_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- Expected accumulative reward when starting from state  $s$  and acting based on  $\pi$

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$

**Output:**  $\pi(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $\pi(s)$  randomly;

**repeat**

    Evaluate  $V_\pi(s), \forall s$ ;

    Improve  $\pi$  such that  $V_\pi(s), \forall s$ , becomes higher;

**until**  $\pi(s)$ 's converge;

**Algorithm 4:** Policy iteration.

# Evaluating $V_\pi(s)$

- How to evaluate the value function of a given  $\pi$ ?

$$V_\pi(s) = \mathbb{E}_{\mathbf{s}^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = s; \pi \right)$$

# Evaluating $V_\pi(s)$

- How to evaluate the value function of a given  $\pi$ ?

$$V_\pi(s) = \mathbb{E}_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- *Bellman expectation equation*:

$$V_\pi(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')], \forall s$$

## M1 Solve the system of linear equations

- $|\mathbb{S}|$  linear equations and  $|\mathbb{S}|$  variables
- Time complexity:  $O(|\mathbb{S}|^3)$

# Evaluating $V_\pi(s)$

- How to evaluate the value function of a given  $\pi$ ?

$$V_\pi(s) = \mathbb{E}_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) | s^{(0)} = s; \pi \right)$$

- *Bellman expectation equation*:

$$V_\pi(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')], \forall s$$

**M1** Solve the system of linear equations

- $|\mathbb{S}|$  linear equations and  $|\mathbb{S}|$  variables
- Time complexity:  $O(|\mathbb{S}|^3)$

**M2** Dynamic programming (just like value iteration):

- Initializes  $V_\pi(s) = 0, \forall s$
- $V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')], \forall s$

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?
- Update rule: for all  $s$  do

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_{\pi}(s')]$$

- Why?

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?
- Update rule: for all  $s$  do

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V_{\pi}(s')]$$

- Why?

$$V_{\pi}(s) = \sum_{s'} P(s'|s; \pi(s))[R(s, \pi(s), s') + \gamma V_{\pi}(s')]$$

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?
- Update rule: for all  $s$  do

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V_{\pi}(s')]$$

- Why?

$$\begin{aligned} V_{\pi}(s) &= \sum_{s'} P(s'|s; \pi(s))[R(s, \pi(s), s') + \gamma V_{\pi}(s')] \\ &\leq \sum_{s'} P(s'|s; \hat{\pi}(s))[R(s, \hat{\pi}(s), s') + \gamma V_{\pi}(s')] \end{aligned}$$

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?
- Update rule: for all  $s$  do

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_{\pi}(s')]$$

- Why?

$$\begin{aligned} V_{\pi}(s) &= \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_{\pi}(s')] \\ &\leq \sum_{s'} P(s'|s; \hat{\pi}(s)) [R(s, \hat{\pi}(s), s') + \gamma V_{\pi}(s')] \\ &\leq \sum_{s'} P(s'|s; \hat{\pi}(s)) \{ R(s, \hat{\pi}(s), s') + \\ &\quad \gamma \sum_{s''} P(s''|s'; \hat{\pi}(s')) [R(s', \hat{\pi}(s'), s'') + \gamma V_{\pi}(s'')] \} \end{aligned}$$

# Policy Improvement

- How to find  $\hat{\pi}$  such that  $V_{\hat{\pi}}(s) \geq V_{\pi}(s)$  for all  $s$ ?
- Update rule: for all  $s$  do

$$\hat{\pi}(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V_{\pi}(s')]$$

- Why?

$$\begin{aligned} V_{\pi}(s) &= \sum_{s'} P(s'|s; \pi(s))[R(s, \pi(s), s') + \gamma V_{\pi}(s')] \\ &\leq \sum_{s'} P(s'|s; \hat{\pi}(s))[R(s, \hat{\pi}(s), s') + \gamma V_{\pi}(s')] \\ &\leq \sum_{s'} P(s'|s; \hat{\pi}(s)) \{ R(s, \hat{\pi}(s), s') + \\ &\quad \gamma \sum_{s''} P(s''|s'; \hat{\pi}(s')) [R(s', \hat{\pi}(s'), s'') + \gamma V_{\pi}(s'')] \} \\ &\dots \\ &\leq E_{s', s'', \dots} (R(s, \hat{\pi}(s), s') + \gamma R(s', \hat{\pi}(s'), s'') + \dots | s^{(0)} = s; \hat{\pi}) \\ &= V_{\hat{\pi}}(s), \forall s \end{aligned}$$

# Algorithm: Policy Iteration

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$

**Output:**  $\pi(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $\pi(s)$  randomly;

**repeat**

    For each state  $s$ , initialize  $V_\pi(s) \leftarrow 0$ ;

**repeat**

**foreach**  $s$  **do**

$| V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')];$

**end**

**until**  $V_\pi(s)$ 's converge;

**foreach**  $s$  **do**

$| \pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')];$

**end**

**until**  $\pi(s)$ 's converge;

**Algorithm 5:** Policy iteration.

# Convergence

## Theorem

*Policy iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

# Convergence

## Theorem

*Policy iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

- Convergence: in every step the policy improves

# Convergence

## Theorem

*Policy iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .*

- Convergence: in every step the policy improves
- Optimal policy: at convergence,  $V_\pi(s)$ ,  $\forall s$ , satisfies the Bellman optimality equation

$$V^*(s) = \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')]$$

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Difference from MDP

- In practice, we may not be able to model the environment as an MDP

# Difference from MDP

- In practice, we may not be able to model the environment as an MDP
- Unknown transition distribution  $P(s'|s; a)$



Kohl and Stone, 2004



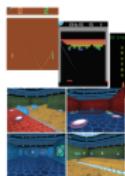
Ng et al, 2004



Tedrake et al, 2005

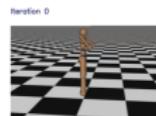


Kober and Peters, 2009



Mnih et al 2013 (DQN)  
Mnih et al, 2015 (A3C)

Silver et al, 2014 (DPG)  
Lillicrap et al, 2015 (DDPG)



Schulman et al,  
2016 (TRPO + GAE)



Levine\*, Finn\*, et  
al, 2016  
(GPS)



Silver\*, Huang\*, et  
al, 2016  
(AlphaGo)

# Difference from MDP

- In practice, we may not be able to model the environment as an MDP
- Unknown transition distribution  $P(s'|s; a)$



Kohl and Stone, 2004



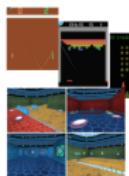
Ng et al, 2004



Tedrake et al, 2005



Kober and Peters, 2009



Mnih et al 2013 (DQN)  
Mnih et al, 2015 (A3C)



Silver et al, 2014 (DPG)  
Lillicrap et al, 2015 (DDPG)



Schulman et al,  
2016 (TRPO + GAE)



Levine\*, Finn\*, et  
al, 2016  
(GPS)



Silver\*, Huang\*, et  
al, 2016  
(AlphaGo)

- Unknown reward function  $R(s, a, s')$



# Exploration vs. Exploitation

- What would you do to get some rewards?



# Exploration vs. Exploitation

- What would you do to get some rewards?
- ① Perform actions randomly to *explore*  $P(s'|s;a)$  and  $R(s,a,s')$  first
  - Collect samples so to estimate  $P(s'|s;a)$  and  $R(s,a,s')$



# Exploration vs. Exploitation

- What would you do to get some rewards?
- ① Perform actions randomly to *explore*  $P(s'|s;a)$  and  $R(s,a,s')$  first
    - Collect samples so to estimate  $P(s'|s;a)$  and  $R(s,a,s')$
  - ② Then, perform actions to *exploit* the learned  $P(s'|s;a)$  and  $R(s,a,s')$ 
    - $\pi^*$  can be computed/planned “in mind” using value/policy iteration



# Model-based RL using Monte Carlo Estimation

- ① Use some *exploration policy*  $\pi'$  to perform one or more *episodes/trails*
  - Each episode records samples of  $P(s'|s;a)$  and  $R(s,a,s')$  from start to terminal state

$$s^{(0)} \xrightarrow{\pi'(s^{(0)})} s^{(1)} \xrightarrow{\pi'(s^{(1)})} \dots \xrightarrow{\pi'(s^{(H-1)})} s^{(H)}$$

and

$$R(s^{(0)}, \pi'(s^{(0)}), s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, \pi(s^{(H-1)}), s^{(H)})$$

- ② Estimate  $P(s'|s;a)$  and  $R(s,a,s')$  using the samples and update the *exploitation policy*  $\pi$ 
  - $\hat{P}(s'|s;a) = \frac{\text{\# times the action } a \text{ takes state } s \text{ to state } s'}{\text{\# times action } a \text{ is taken in state } s}$
  - $\hat{R}(s,a,s') = \text{average of reward values received when } a \text{ takes } s \text{ to } s'$
- ③ Repeat from Step 1, but gradually mix  $\pi$  into  $\pi'$ 
  - Mix-in strategy? E.g.,  $\epsilon$ -greedy (more on this later)

# Problems of Model-based RL

- There may be lots of  $P(s'|s; \mathbf{a})$  and  $R(s, \mathbf{a}, s')$  to estimate
- If  $P(s'|s; \mathbf{a})$  is low, there may be too few samples to have a good estimate
- Low  $P(s'|s; \mathbf{a})$  may also lead to a poor estimate of  $R(s, \mathbf{a}, s')$ 
  - when  $R(s, \mathbf{a}, s')$  depends on  $s'$

# Problems of Model-based RL

- There may be lots of  $P(s'|s; \mathbf{a})$  and  $R(s, \mathbf{a}, s')$  to estimate
- If  $P(s'|s; \mathbf{a})$  is low, there may be too few samples to have a good estimate
- Low  $P(s'|s; \mathbf{a})$  may also lead to a poor estimate of  $R(s, \mathbf{a}, s')$ 
  - when  $R(s, \mathbf{a}, s')$  depends on  $s'$
- We estimate  $P(s'|s; \mathbf{a})$ 's and  $R(s, \mathbf{a}, s')$ 's in order to compute  $V^*(s)/V_\pi(s)$  and solve  $\pi^*$
- Why not estimate  $V^*(s)/V_\pi(s)$  directly?

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Model-based vs. Model-Free

- How to estimate  $E(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})f(\mathbf{x})$  given samples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots$ ?

# Model-based vs. Model-Free

- How to estimate  $E(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})f(\mathbf{x})$  given samples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots$ ?
- **Model-based** estimation:
  - ① For each  $\mathbf{x}$ , estimate  $\hat{P}(\mathbf{x}|\mathbf{y}) = \frac{\text{count}(\mathbf{x}^{(i)}=\mathbf{x} \text{ and } \mathbf{y}^{(i)}=\mathbf{y})}{\text{count}(\mathbf{y}^{(j)}=\mathbf{y})}$
  - ②  $\hat{E}(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} \hat{P}(\mathbf{x}|\mathbf{y})f(\mathbf{x})$

# Model-based vs. Model-Free

- How to estimate  $E(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})f(\mathbf{x})$  given samples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots$ ?
- **Model-based** estimation:
  - ① For each  $\mathbf{x}$ , estimate  $\hat{P}(\mathbf{x}|\mathbf{y}) = \frac{\text{count}(\mathbf{x}^{(i)} = \mathbf{x} \text{ and } \mathbf{y}^{(i)} = \mathbf{y})}{\text{count}(\mathbf{y}^{(j)} = \mathbf{y})}$
  - ②  $\hat{E}(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} \hat{P}(\mathbf{x}|\mathbf{y})f(\mathbf{x})$
- **Model-free** estimation:

$$\hat{E}(f(\mathbf{x})|\mathbf{y}) = \frac{1}{\text{count}(\mathbf{y}^{(j)} = \mathbf{y})} \sum_{i: \mathbf{x}^{(i)} = \mathbf{x} \text{ and } \mathbf{y}^{(i)} = \mathbf{y}} f(\mathbf{x}^{(i)})$$

- Why does it work?

# Model-based vs. Model-Free

- How to estimate  $E(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})f(\mathbf{x})$  given samples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots$ ?
- **Model-based** estimation:

- ① For each  $\mathbf{x}$ , estimate  $\hat{P}(\mathbf{x}|\mathbf{y}) = \frac{\text{count}(\mathbf{x}^{(i)} = \mathbf{x} \text{ and } \mathbf{y}^{(i)} = \mathbf{y})}{\text{count}(\mathbf{y}^{(j)} = \mathbf{y})}$
- ②  $\hat{E}(f(\mathbf{x})|\mathbf{y}) = \sum_{\mathbf{x}} \hat{P}(\mathbf{x}|\mathbf{y})f(\mathbf{x})$

- **Model-free** estimation:

$$\hat{E}(f(\mathbf{x})|\mathbf{y}) = \frac{1}{\text{count}(\mathbf{y}^{(j)} = \mathbf{y})} \sum_{i: \mathbf{x}^{(i)} = \mathbf{x} \text{ and } \mathbf{y}^{(i)} = \mathbf{y}} f(\mathbf{x}^{(i)})$$

- Why does it work? Because samples appear in right frequencies

# Challenges of Model-Free RL

- Value iteration: start from  $V^*(s) \leftarrow 0$ 
  - ① Iterate  $V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$  until converge
  - ② Solve  $\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$
- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ① Iterate  $V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s))[R(s, \pi(s), s') + \gamma V_\pi(s')]$  until converge
  - ② Solve  $\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V_\pi(s')]$

# Challenges of Model-Free RL

- Value iteration: start from  $V^*(s) \leftarrow 0$ 
  - ① Iterate  $V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$  until converge
    - ⌚ **Not easy to estimate**  $V^*(s) = \max_{\pi} E\left(\sum_t \gamma^t R^{(t)} | s^{(0)} = s; \pi\right)$
  - ② Solve  $\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$ 
    - ⌚ **Need model to solve**
- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ① Iterate  $V_{\pi}(s) \leftarrow \sum_{s'} P(s'|s; \pi(s))[R(s, \pi(s), s') + \gamma V_{\pi}(s')]$  until converge
    - ⌚ **Can estimate**  $V_{\pi}(s) = E\left(\sum_t \gamma^t R^{(t)} | s^{(0)} = s; \pi\right)$  **using MC est.**
  - ② Solve  $\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V_{\pi}(s')]$ 
    - ⌚ **Need model to solve**

# $Q$ Function for $\pi$

- Value function for a given  $\pi$ :

$$V_\pi(s) = \mathbb{E}_{\mathbf{s}^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = s; \pi \right)$$

with recurrence (Bellman expectation equation):

$$V_\pi(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')], \forall s$$

- Maximum expected accumulative reward when starting from state  $s$  and acting based on  $\pi$  onward

# $Q$ Function for $\pi$

- Value function for a given  $\pi$ :

$$V_\pi(s) = \mathbb{E}_{\mathbf{s}^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = s; \pi \right)$$

with recurrence (Bellman expectation equation):

$$V_\pi(s) = \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')], \forall s$$

- Maximum expected accumulative reward when starting from state  $s$  and acting based on  $\pi$  onward
- Define  **$Q$  function** for  $\pi$ :

$$Q_\pi(s, a) = \mathbb{E}_{\mathbf{s}^{(1)}, \dots} \left( R(s, a, \mathbf{s}^{(1)}) + \sum_{t=1}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); s, a, \pi \right)$$

such that  $V_\pi(s) = Q_\pi(s, \pi(s))$  with recurrence:

$$Q_\pi(s, a) = \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma Q_\pi(s', \pi(s'))], \forall s, a$$

- Maximum expected accumulative reward when starting from state  $s$ , taking action  $a$ , and then acting based on  $\pi$

# Algorithm: Policy Iteration based on $Q_\pi$

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$

**Output:**  $\pi(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $\pi(s)$  randomly;

**repeat**

    For each state  $s$ , initialize  $V_\pi(s) \leftarrow 0$  Initialize  $Q_\pi(s, a) = 0, \forall s, a$ ;

**repeat**

**foreach**  $s$  **and**  $a$  **do**

$$V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')]$$

$$Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma Q_\pi(s', \pi(s'))];$$

**end**

**until**  $V_\pi(s)$ 's  $Q_\pi(s, a)$ 's converge;

**foreach**  $s$  **do**

$$\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')]$$

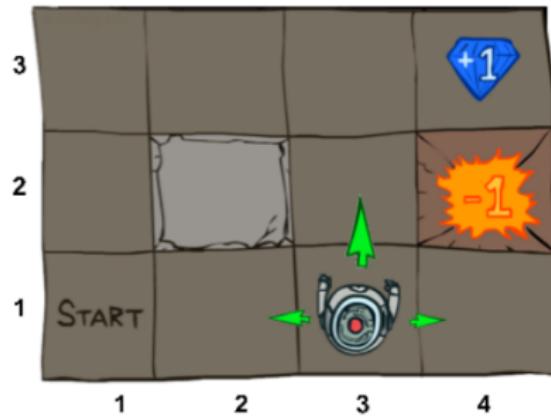
$$\pi(s) \leftarrow \arg \max_a Q_\pi(s, a);$$

**end**

**until**  $\pi(s)$ 's converge;

**Algorithm 6:** Policy iteration.

# Example



0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.43	0.48	0.28

VALUES AFTER 100 ITERATIONS

0.59	0.67	0.77	1.00
0.57	0.64	0.60	0.74
0.53	0.57	0.67	0.66
0.57	0.51	0.53	0.85
0.51	0.51	-0.60	-1.00
0.46	0.49	0.30	
0.49	0.41	0.40	0.48
0.45	0.43	0.42	0.29
0.44	0.40	0.40	0.28
0.41	0.41	0.40	0.13

Q-VALUES AFTER 100 ITERATIONS

# Policy Iteration and Model-Free RL

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_{\mathbf{a}} Q_\pi(s, \mathbf{a})$

# Policy Iteration and Model-Free RL

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
    - ⌚ **Can estimate**  $Q_\pi(s, a) = E\left(R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)}\right)$  **using MC est.**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$ 
    - ⌚ **No need for model to solve**

# Policy Iteration and Model-Free RL

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
    - ⌚ **Can estimate**  $Q_\pi(s, a) = E\left(R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)}\right)$  **using MC est.**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$ 
    - ⌚ **No need for model to solve**
- RL: start from a random  $\pi$  **for both exploration & exploitation**, repeat until converge:
  - ➊ Create episodes using  $\pi$  and get MC estimates  $\hat{Q}_\pi(s, a), \forall s, a$
  - ➋ Update  $\pi$  by  $\pi(s) \leftarrow \arg \max_a \hat{Q}_\pi(s, a), \forall s$

# Policy Iteration and Model-Free RL

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
    - ➋ **Can estimate  $Q_\pi(s, a) = E(R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)})$  using MC est.**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$ 
    - ➋ **No need for model to solve**
- RL: start from a random  $\pi$  **for both exploration & exploitation**, repeat until converge:
  - ➊ Create episodes using  $\pi$  and get MC estimates  $\hat{Q}_\pi(s, a), \forall s, a$
  - ➋ Update  $\pi$  by  $\pi(s) \leftarrow \arg \max_a \hat{Q}_\pi(s, a), \forall s$
- Problem:  $\pi$  improves little after running lots of trials

# Policy Iteration and Model-Free RL

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
    - ⌚ **Can estimate**  $Q_\pi(s, a) = E\left(R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)}\right)$  **using MC est.**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$ 
    - ⌚ **No need for model to solve**
- RL: start from a random  $\pi$  **for both exploration & exploitation**, repeat until converge:
  - ➊ Create episodes using  $\pi$  and get MC estimates  $\hat{Q}_\pi(s, a), \forall s, a$
  - ➋ Update  $\pi$  by  $\pi(s) \leftarrow \arg \max_a \hat{Q}_\pi(s, a), \forall s$
- Problem:  $\pi$  improves little after running lots of trials
- Can we improve  $\pi$  right after each action?

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Policy Iteration Revisited

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_{\mathbf{a}} Q_\pi(s, \mathbf{a})$

# Policy Iteration Revisited

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge  
    ⌚ **Can estimate**  $Q_\pi(s, a) = E \left( R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)} \right)$  **using MC est.**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$   
    ⌚ **No need for model to solve**

# Policy Iteration Revisited

- Policy iteration: start from a random  $\pi$ , repeat until converge:
  - ➊ Iterate  $Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$  until converge
    - ⌚ **Can estimate**  $Q_\pi(s, a) = E\left(R^{(0)} + \sum_{t=1}^{\infty} \gamma^t R^{(t)}\right)$  **using MC est.**
    - ⌚ **Can also estimate**  $Q_\pi(s, a)$  **based on the recurrence**
  - ➋ Solve  $\pi(s) \leftarrow \arg \max_a Q_\pi(s, a)$ 
    - ⌚ **No need for model to solve**

# Temporal Difference Estimation I

$$Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$$

- Given the samples

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, a^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q_\pi(s, a)$ :

①  $\hat{Q}_\pi(s, a) \leftarrow$  random value,  $\forall s, a$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\begin{aligned} \hat{Q}_\pi(s^{(t)}, a^{(t)}) &\leftarrow \hat{Q}_\pi(s^{(t)}, a^{(t)}) + \\ &\eta \left[ (R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \hat{Q}_\pi(s^{(t+1)}, \pi(s^{(t+1)}))) - \hat{Q}_\pi(s^{(t)}, a^{(t)}) \right] \end{aligned}$$

# Temporal Difference Estimation I

$$Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$$

- Given the samples

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, a^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q_\pi(s, a)$ :

①  $\hat{Q}_\pi(s, a) \leftarrow$  random value,  $\forall s, a$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\begin{aligned} \hat{Q}_\pi(s^{(t)}, a^{(t)}) &\leftarrow \hat{Q}_\pi(s^{(t)}, a^{(t)}) + \\ &\eta \left[ (R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \hat{Q}_\pi(s^{(t+1)}, \pi(s^{(t+1)}))) - \hat{Q}_\pi(s^{(t)}, a^{(t)}) \right] \end{aligned}$$

- $\hat{Q}_\pi(s, a)$  can be updated **on the fly** during an episode

# Temporal Difference Estimation I

$$Q_\pi(s, a) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma Q_\pi(s', \pi(s'))]$$

- Given the samples

$$s^{(0)} \xrightarrow{a^{(0)}} s^{(1)} \xrightarrow{a^{(1)}} \dots \xrightarrow{a^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, a^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, a^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q_\pi(s, a)$ :

①  $\hat{Q}_\pi(s, a) \leftarrow$  random value,  $\forall s, a$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\hat{Q}_\pi(s^{(t)}, a^{(t)}) \leftarrow \hat{Q}_\pi(s^{(t)}, a^{(t)}) + \eta \left[ (R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \hat{Q}_\pi(s^{(t+1)}, \pi(s^{(t+1)}))) - \hat{Q}_\pi(s^{(t)}, a^{(t)}) \right]$$

- $\hat{Q}_\pi(s, a)$  can be updated **on the fly** during an episode
- $\eta$  the “learning rate”?

# Temporal Difference Estimation II

$$\begin{aligned}\hat{Q}_\pi(s, a) &\leftarrow \hat{Q}_\pi(s, a) + \eta \left[ (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) - \hat{Q}_\pi(s, a) \right], \\ &= \eta (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) + (1 - \eta) \hat{Q}_\pi(s, a)\end{aligned}$$

# Temporal Difference Estimation II

$$\begin{aligned}\hat{Q}_\pi(s, a) &\leftarrow \hat{Q}_\pi(s, a) + \eta \left[ (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) - \hat{Q}_\pi(s, a) \right], \\ &= \eta (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) + (1 - \eta) \hat{Q}_\pi(s, a)\end{aligned},$$

- **Exponential moving average:**

$$\bar{x}_n = \frac{x^{(n)} + (1 - \eta)x^{(n-1)} + (1 - \eta)^2x^{(n-2)} + \dots}{1 + (1 - \eta) + (1 - \eta)^2 + \dots},$$

where  $\eta \in [0, 1]$  is the “forget rate”

- Recent samples are exponentially more important

# Temporal Difference Estimation II

$$\begin{aligned}\hat{Q}_\pi(s, a) &\leftarrow \hat{Q}_\pi(s, a) + \eta \left[ (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) - \hat{Q}_\pi(s, a) \right], \\ &= \eta (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) + (1 - \eta) \hat{Q}_\pi(s, a)\end{aligned},$$

- **Exponential moving average:**

$$\bar{x}_n = \frac{x^{(n)} + (1 - \eta)x^{(n-1)} + (1 - \eta)^2x^{(n-2)} + \dots}{1 + (1 - \eta) + (1 - \eta)^2 + \dots},$$

where  $\eta \in [0, 1]$  is the “forget rate”

- Recent samples are exponentially more important
- Since  $1/\eta = 1 + (1 - \eta) + (1 - \eta)^2 + \dots$ , we have  
 $\bar{x}_n = \eta x^{(n)} + (1 - \eta)\bar{x}_{n-1}$  [Proof]

# Temporal Difference Estimation II

$$\begin{aligned}\hat{Q}_\pi(s, a) &\leftarrow \hat{Q}_\pi(s, a) + \eta \left[ (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) - \hat{Q}_\pi(s, a) \right], \\ &= \eta (R(s, a, s') + \gamma \hat{Q}_\pi(s', \pi(s'))) + (1 - \eta) \hat{Q}_\pi(s, a)\end{aligned},$$

- **Exponential moving average:**

$$\bar{x}_n = \frac{x^{(n)} + (1 - \eta)x^{(n-1)} + (1 - \eta)^2x^{(n-2)} + \dots}{1 + (1 - \eta) + (1 - \eta)^2 + \dots},$$

where  $\eta \in [0, 1]$  is the “forget rate”

- Recent samples are exponentially more important
- Since  $1/\eta = 1 + (1 - \eta) + (1 - \eta)^2 + \dots$ , we have  
 $\bar{x}_n = \eta x^{(n)} + (1 - \eta)\bar{x}_{n-1}$  [Proof]
- As long as  $\eta$  gradually decreases to 0:
  - $\hat{Q}_\pi(s, a)$  degenerates to average of accumulative rewards
  - $\hat{Q}_\pi(s, a)$  converges to  $Q_\pi(s, a)$  (more on this later)

# Algorithm: SARSA (Simplified)

**Input:**  $\mathbb{S}$ ,  $\mathbb{A}$ , and  $\gamma$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$  and  $a$ , initialize  $Q_\pi(s, a)$  arbitrarily;

**foreach** episode **do**

    Set  $s$  to initial state;

**repeat**

        Take action  $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$ ;

        Observe  $s'$  and reward  $R(s, a, s')$ ;

$Q_\pi(s, a) \leftarrow$

$Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s')) - Q_\pi(s, a))]$ ;

$s \leftarrow s'$ ;

**until**  $s$  is terminal state;

**end**

**Algorithm 7:** State-Action-Reward-State-Action (SARSA).

# Algorithm: SARSA (Simplified)

**Input:**  $\mathbb{S}$ ,  $\mathbb{A}$ , and  $\gamma$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$  and  $a$ , initialize  $Q_\pi(s, a)$  arbitrarily;

**foreach** episode **do**

    Set  $s$  to initial state;

**repeat**

        Take action  $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$ ;

        Observe  $s'$  and reward  $R(s, a, s')$ ;

$Q_\pi(s, a) \leftarrow$

$Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s')) - Q_\pi(s, a))]$ ;

$s \leftarrow s'$ ;

**until**  $s$  is terminal state;

**end**

**Algorithm 8:** State-Action-Reward-State-Action (SARSA).

- Policy improves **each time** when deciding the next action  $a$

# Convergence

## Theorem

SARSA converges and gives the optimal policy  $\pi^*$  almost surely if

- 1)  $\pi$  is GLIE (Greedy in the Limit with Infinite Exploration);
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

# Convergence

## Theorem

SARSA converges and gives the optimal policy  $\pi^*$  almost surely if

- 1)  $\pi$  is GLIE (Greedy in the Limit with Infinite Exploration);
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

- Greedy in the limit: the policy  $\pi$  converges (in the limit) to the exploitation/greedy policy
  - At each step, we choose  $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$

# Convergence

## Theorem

SARSA converges and gives the optimal policy  $\pi^*$  almost surely if

- 1)  $\pi$  is GLIE (Greedy in the Limit with Infinite Exploration);
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

- Greedy in the limit: the policy  $\pi$  converges (in the limit) to the exploitation/greedy policy
  - At each step, we choose  $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$
- Infinite exploration: all  $(s, a)$  pairs are visited infinite times
  - $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$  **cannot** guarantee this
  - Need a better way (next section)

# Convergence

## Theorem

SARSA converges and gives the optimal policy  $\pi^*$  almost surely if

- 1)  $\pi$  is GLIE (Greedy in the Limit with Infinite Exploration);
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

- Greedy in the limit: the policy  $\pi$  converges (in the limit) to the exploitation/greedy policy
  - At each step, we choose  $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$
- Infinite exploration: all  $(s, a)$  pairs are visited infinite times
  - $a \leftarrow \arg \max_{a'} Q_\pi(s, a')$  **cannot** guarantee this
  - Need a better way (next section)
- Furthermore,  $\eta$  should satisfy:  $\sum_t \eta^{(t)} = \infty$  and  $\sum_t \eta^{(t)2} < \infty$ 
  - E.g.,  $\eta^{(t)} = O(\frac{1}{t})$
  - $\sum_t \frac{1}{t}$  is a [harmonic series](#) known to diverge
  - $\sum_t (\frac{1}{t})^p$ ,  $p > 1$ , is a  $p$ -series converging to [Riemann zeta](#)  $\zeta(p)$

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# General RL Steps

- Repeat until converge:
  - ① Use some exploration policy  $\pi'$  to run episodes/trails
  - ② Estimate targets (e.g.,  $P(s'|s;a)$ ,  $R(s,a,s')$ , or  $Q_\pi(s,a)$ ) and update the exploitation policy  $\pi$
  - ③ Gradually mix  $\pi$  into  $\pi'$

# General RL Steps

- Repeat until converge:
  - ① Use some exploration policy  $\pi'$  to run episodes/trails
  - ② Estimate targets (e.g.,  $P(s'|s;a)$ ,  $R(s,a,s')$ , or  $Q_\pi(s,a)$ ) and update the exploitation policy  $\pi$
  - ③ Gradually mix  $\pi$  into  $\pi'$
- Goals of mix-in/exploration strategy:
  - Infinite exploration with  $\pi'$
  - $\pi'$  is greedy/exploitative in the end

# $\varepsilon$ -Greedy Strategy

- At every time step, flip a coin
  - With probability  $\varepsilon$ , act randomly (explore)
  - With probability  $(1 - \varepsilon)$ , compute/update the exploitation policy and act accordingly (exploit)

# $\varepsilon$ -Greedy Strategy

- At every time step, flip a coin
  - With probability  $\varepsilon$ , act randomly (explore)
  - With probability  $(1 - \varepsilon)$ , compute/update the exploitation policy and act accordingly (exploit)
- Gradually decrease  $\varepsilon$  over time

# $\varepsilon$ -Greedy Strategy

- At every time step, flip a coin
  - With probability  $\varepsilon$ , act randomly (explore)
  - With probability  $(1 - \varepsilon)$ , compute/update the exploitation policy and act accordingly (exploit)
- Gradually decrease  $\varepsilon$  over time
- At each time step, the action is at either of two extremes
  - Exploration or exploitation
- “Soft” policy between the two extremes?

# Softmax Strategy

- Idea: perform action  $a$  more often if  $a$  gives more accumulative rewards
- E.g., in SARSA, choose  $a$  from  $s$  by sampling from the distribution:

$$P(a|s) = \frac{\exp(Q_\pi(s, a)/t)}{\sum_{a'}(\exp Q_\pi(s, a')/t)}, \forall a$$

- Softmax function converts  $Q_\pi(s, a)$ 's to probabilities

# Softmax Strategy

- Idea: perform action  $a$  more often if  $a$  gives more accumulative rewards
- E.g., in SARSA, choose  $a$  from  $s$  by sampling from the distribution:

$$P(a|s) = \frac{\exp(Q_\pi(s, a)/t)}{\sum_{a'}(\exp Q_\pi(s, a')/t)}, \forall a$$

- Softmax function converts  $Q_\pi(s, a)$ 's to probabilities
- **Temperature**  $t$  starts from a high value (exploration), and decreases over time (exploitation)

# Exploration Function

- Idea: to explore areas with fewest samples
- E.g., in each step of SARSA, define an exploration function

$$f(q, n) = q + K/n,$$

where

- $q$  the estimated  $Q$ -value
- $n$  the number of samples for the estimate
- $K$  some positive constant

# Exploration Function

- Idea: to explore areas with fewest samples
- E.g., in each step of SARSA, define an exploration function

$$f(q, n) = q + K/n,$$

where

- $q$  the estimated  $Q$ -value
- $n$  the number of samples for the estimate
- $K$  some positive constant
- Instead of:

$$Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', a')) - Q_\pi(s, a)]$$

- Use  $f$  when updating  $Q_\pi(s, a)$ :

$$Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta \{ [R(s, a, s') + \gamma f(Q_\pi(s', a'), \text{count}(s', a'))] - Q_\pi(s, a) \}$$

- Infinite exploration
- Exploit once exploring enough

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- *Q*-Learning (Model-Free)
- SARSA vs. *Q*-Learning

# Model-Free RL with Value Iteration?

- Value iteration: start from  $V^*(s) \leftarrow 0$ 
  - ① Iterate  $V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$  until converge  
    ⌚ **Not easy to estimate**  $V^*(s) = \max_\pi E\left(\sum_t \gamma^t R^{(t)} | s^{(0)} = s; \pi\right)$
  - ② Solve  $\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$   
    ⌚ **Need model to solve**

# Model-Free RL with Value Iteration?

- Value iteration: start from  $V^*(s) \leftarrow 0$ 
  - ➊ Iterate  $V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')] \text{ until converge}$   
➋ **Not easy to estimate**  $V^*(s) = \max_\pi E\left(\sum_t \gamma^t R^{(t)} | s^{(0)} = s; \pi\right)$
  - ➋ Solve  $\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma V^*(s')]$   
**Need model to solve**
- $Q$ -version that helps sample-based estimation?

# Optimal $Q$ Function

- Optimal value function:

$$V^*(s) = \max_{\pi} E_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) \mid s^{(0)} = s; \pi \right)$$

with recurrence

$$V^*(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^*(s')], \forall s$$

- Maximum expected accumulative reward when starting from state  $s$  and acting optimally onward

# Optimal $Q$ Function

- Optimal value function:

$$V^*(s) = \max_{\pi} E_{s^{(1)}, \dots} \left( \sum_{t=0}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) \mid s^{(0)} = s; \pi \right)$$

with recurrence

$$V^*(s) = \max_{\mathbf{a}} \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma V^*(s')], \forall s$$

- Maximum expected accumulative reward when starting from state  $s$  and acting optimally onward
- $Q^*$  function:**

$$Q^*(s, \mathbf{a}) = \max_{\pi} E_{s^{(1)}, \dots} \left( R(s, \mathbf{a}, s^{(1)}) + \sum_{t=1}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) ; s, \mathbf{a}, \pi \right)$$

such that  $V^*(s) = \max_{\mathbf{a}} Q^*(s, \mathbf{a})$  with recurrence:

$$Q^*(s, \mathbf{a}) = \sum_{s'} P(s' | s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')], \forall s$$

- Maximum expected accumulative reward when starting from state  $s$ , taking action  $\mathbf{a}$ , and then acting optimally onward

# Algorithm: $Q$ -Value Iteration

**Input:** MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$ , initialize  $V^*(s) \leftarrow 0$  Initialize  $Q^*(s, a) = 0, \forall s, a$ ;

repeat

    foreach  $s$  and  $a$  do

$$V^*(s) \leftarrow \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')]$$

$$Q^*(s, a) \leftarrow \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')];$$

    end

until  $V^*(s)$ 's  $Q^*(s, a)$ 's converge;

foreach  $s$  do

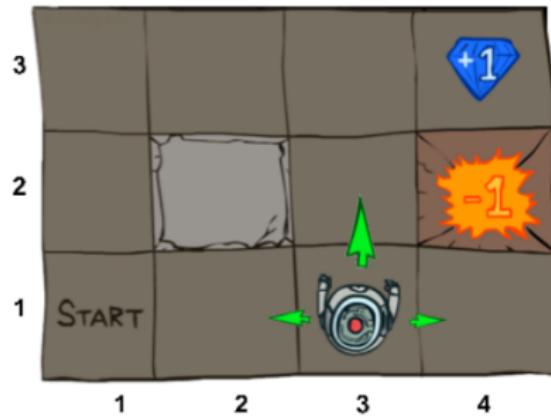
$$\pi^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V^*(s')]$$

$$\pi^*(s) \leftarrow \arg \max_a Q^*(s, a);$$

end

**Algorithm 9:**  $Q$ -Value iteration with infinite horizon.

# Example



0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.43	0.48	0.28

VALUES AFTER 100 ITERATIONS

0.59	0.67	0.77	1.00
0.57	0.64	0.60	0.74
0.53	0.57	0.67	0.66
0.51	0.51		0.85
0.46	0.49	0.57	
0.45	0.41	0.53	-0.60
0.44	0.43	0.30	
0.40	0.42	0.48	-0.65
0.40	0.40	0.29	0.28
0.41	0.42	0.41	0.13

Q-VALUES AFTER 100 ITERATIONS

# $Q$ -Value Iteration

- Value iteration: start from  $Q^*(s, \mathbf{a}) \leftarrow 0$ 
  - ➊ Iterate  $Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a})[R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')] \text{ until converge}$
  - ➋ Solve  $\pi^*(s) \leftarrow \arg \max_{\mathbf{a}} Q^*(s, \mathbf{a})$

# $Q$ -Value Iteration

- Value iteration: start from  $Q^*(s, \mathbf{a}) \leftarrow 0$ 
  - ① Iterate  $Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a})[R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')]$  until converge
    - ⌚ **Still not easy to estimate**  $Q^*(s, \mathbf{a}) = \max_{\pi} E \left( R^{(0)} + \sum_t \gamma^t R^{(t)} \right)$
  - ② Solve  $\pi^*(s) \leftarrow \arg \max_{\mathbf{a}} Q^*(s, \mathbf{a})$ 
    - ⌚ **No need for model to solve**

# $Q$ -Value Iteration

- Value iteration: start from  $Q^*(s, a) \leftarrow 0$ 
  - ① Iterate  $Q^*(s, a) \leftarrow \sum_{s'} P(s'|s; a)[R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$  until converge
    - ⌚ Still not easy to estimate  $Q^*(s, a) = \max_{\pi} E(R^{(0)} + \sum_t \gamma^t R^{(t)})$
    - ⌚ But we can estimate  $Q^*(s, a)$  based on the recurrence now!
  - ② Solve  $\pi^*(s) \leftarrow \arg \max_a Q^*(s, a)$ 
    - ⌚ No need for model to solve

# Temporal Difference Estimation

$$Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')]$$

- Given an exploration policy  $\pi'$  and samples

$$s^{(0)} \xrightarrow{\mathbf{a}^{(0)}} s^{(1)} \xrightarrow{\mathbf{a}^{(1)}} \dots \xrightarrow{\mathbf{a}^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, \mathbf{a}^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, \mathbf{a}^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q^*(s, \mathbf{a})$  for exploitation policy  $\pi$ :

①  $\hat{Q}^*(s, \mathbf{a}) \leftarrow$  random value,  $\forall s, \mathbf{a}$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\begin{aligned}\hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) &\leftarrow \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) + \\ &\quad \eta \left[ (R(s^{(t)}, \mathbf{a}^{(t)}, s^{(t+1)}) + \gamma \max_{\mathbf{a}} \hat{Q}^*(s^{(t+1)}, \mathbf{a})) - \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) \right]\end{aligned}$$

# Temporal Difference Estimation

$$Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')]$$

- Given an exploration policy  $\pi'$  and samples

$$s^{(0)} \xrightarrow{\mathbf{a}^{(0)}} s^{(1)} \xrightarrow{\mathbf{a}^{(1)}} \dots \xrightarrow{\mathbf{a}^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, \mathbf{a}^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, \mathbf{a}^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q^*(s, \mathbf{a})$  for exploitation policy  $\pi$ :

①  $\hat{Q}^*(s, \mathbf{a}) \leftarrow$  random value,  $\forall s, \mathbf{a}$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\begin{aligned}\hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) &\leftarrow \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) + \\ \eta \left[ (R(s^{(t)}, \mathbf{a}^{(t)}, s^{(t+1)}) + \gamma \max_{\mathbf{a}} \hat{Q}^*(s^{(t+1)}, \mathbf{a})) - \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) \right]\end{aligned}$$

- $\hat{Q}^*(s, \mathbf{a})$  can be updated **on the fly** during exploration

# Temporal Difference Estimation

$$Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} Q^*(s', \mathbf{a}')]$$

- Given an exploration policy  $\pi'$  and samples

$$s^{(0)} \xrightarrow{\mathbf{a}^{(0)}} s^{(1)} \xrightarrow{\mathbf{a}^{(1)}} \dots \xrightarrow{\mathbf{a}^{(H-1)}} s^{(H)}$$

and

$$R(s^{(0)}, \mathbf{a}^{(0)}, s^{(1)}) \rightarrow \dots \rightarrow R(s^{(H-1)}, \mathbf{a}^{(H-1)}, s^{(H)})$$

- Temporal difference (TD) estimation** of  $Q^*(s, \mathbf{a})$  for exploitation policy  $\pi$ :

①  $\hat{Q}^*(s, \mathbf{a}) \leftarrow$  random value,  $\forall s, \mathbf{a}$

② Repeat until converge **for each action  $a^{(t)}$** :

$$\hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) \leftarrow \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) + \eta \left[ (R(s^{(t)}, \mathbf{a}^{(t)}, s^{(t+1)}) + \gamma \max_{\mathbf{a}} \hat{Q}^*(s^{(t+1)}, \mathbf{a})) - \hat{Q}^*(s^{(t)}, \mathbf{a}^{(t)}) \right]$$

- $\hat{Q}^*(s, \mathbf{a})$  can be updated **on the fly** during exploration

- $\eta$  the “forget rate” of moving avg. that gradually decreases

# Algorithm: $Q$ -Learning

**Input:**  $\mathbb{S}$ ,  $\mathbb{A}$ , and  $\gamma$

**Output:**  $\pi^*(s)$ 's for all  $s$ 's

For each state  $s$  and  $a$ , initialize  $Q^*(s, a)$  arbitrarily;

**foreach** episode **do**

    Set  $s$  to initial state;

**repeat**

        Take action  $a$  from  $s$  using some exploration policy  $\pi'$   
        derived from  $Q^*$  (e.g.,  $\varepsilon$ -greedy);

        Observe  $s'$  and reward  $R(s, a, s')$ ;

$Q^*(s, a) \leftarrow$

$Q^*(s, a) + \eta [(R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) - Q^*(s, a)]$ ;

$s \leftarrow s'$ ;

**until**  $s$  is terminal state;

**end**

**Algorithm 10:**  $Q$ -learning.

# Convergence

## Theorem

*Q-learning converges and gives the optimal policy  $\pi^*$  if*

- 1)  $\pi'$  has explored enough;
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

# Convergence

## Theorem

*Q-learning converges and gives the optimal policy  $\pi^*$  if*

- 1)  $\pi'$  has explored enough;
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

- $\pi'$  has explored enough:
  - All states and actions are visited infinitely often
  - Does **not** matter how  $\pi'$  selects actions!

# Convergence

## Theorem

*Q-learning converges and gives the optimal policy  $\pi^*$  if*

- 1)  $\pi'$  has explored enough;
- 2)  $\eta$  small enough eventually, but not decreasing too fast.

- $\pi'$  has explored enough:
  - All states and actions are visited infinitely often
  - Does **not** matter how  $\pi'$  selects actions!
- $\eta$  satisfies  $\sum_t \eta^{(t)} = \infty$  and  $\sum_t \eta^{(t)2} < \infty$ 
  - E.g.,  $\eta^{(t)} = O(\frac{1}{t})$

# Outline

## ① Introduction

## ② Markov Decision Process

- Value Iteration
- Policy Iteration

## ③ Reinforcement Learning

- Model-Free RL using Monte Carlo Estimation
- Temporal-Difference Estimation and SARSA (Model-Free)
- Exploration Strategies
- $Q$ -Learning (Model-Free)
- SARSA vs.  $Q$ -Learning

# Off-Policy vs. On-Policy RL

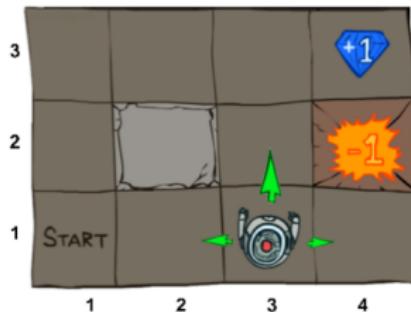
- General RL steps: repeat until converge:
  - 1 Use some exploration policy  $\pi'$  to run episodes/trails
  - 2 Use samples to update the exploitation policy  $\pi$
  - 3 Gradually mix  $\pi$  into  $\pi'$

# Off-Policy vs. On-Policy RL

- General RL steps: repeat until converge:
  - 1 Use some exploration policy  $\pi'$  to run episodes/trails
  - 2 Use samples to update the exploitation policy  $\pi$
  - 3 Gradually mix  $\pi$  into  $\pi'$
- Off-policy:**  $\pi$  updated toward a greedy policy *independent with*  $\pi'$ 
  - $Q$ -learning:  
$$Q^*(s, a) \leftarrow Q^*(s, a) + \eta [(R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) - Q^*(s, a)]$$
- On-policy:**  $\pi$  updated to improve (and *depends on*)  $\pi'$ 
  - SARSA:  $Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s'))) - Q_\pi(s, a)]$

# Off-Policy vs. On-Policy RL

- General RL steps: repeat until converge:
  - Use some exploration policy  $\pi'$  to run episodes/trails
  - Use samples to update the exploitation policy  $\pi$
  - Gradually mix  $\pi$  into  $\pi'$
- Off-policy:**  $\pi$  updated toward a greedy policy *independent with*  $\pi'$ 
  - $Q$ -learning:  
$$Q^*(s, a) \leftarrow Q^*(s, a) + \eta [(R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) - Q^*(s, a)]$$
- On-policy:**  $\pi$  updated to improve (and *depends on*)  $\pi'$ 
  - SARSA:  $Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s'))) - Q_\pi(s, a)]$

 $\pi'$  $\pi$

# Practical Results

- SARSA has the capability to *avoid the mistakes due to exploration*
  - E.g., the [maze-with-cliff problem](#)
  - Advantageous when  $\varepsilon > 0$  with  $\varepsilon$ -greedy exploration strategy
  - Converges to optimal policy  $\pi^*$  (as  $Q$ -learning) when  $\varepsilon \rightarrow 0$

# Practical Results

- SARSA has the capability to *avoid the mistakes due to exploration*
  - E.g., the maze-with-cliff problem
  - Advantageous when  $\varepsilon > 0$  with  $\varepsilon$ -greedy exploration strategy
  - Converges to optimal policy  $\pi^*$  (as  $Q$ -learning) when  $\varepsilon \rightarrow 0$
- But  $Q$ -learning has the capability to *continue learning while changing the exploration policy*

# Remarks

- Update rule for  $Q^*$  (or  $Q_\pi$ ) takes terminal states differently:

$$Q^*(s, a) \leftarrow \begin{cases} (1 - \eta)Q^*(s, a) + \eta R(s, a, s') & \text{if } s' \text{ is terminal} \\ (1 - \eta)Q^*(s, a) + \eta [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')] & \text{otherwise} \end{cases}$$

- Better convergence in practice

# Remarks

- Update rule for  $Q^*$  (or  $Q_\pi$ ) takes terminal states differently:

$$Q^*(s, a) \leftarrow \begin{cases} (1 - \eta)Q^*(s, a) + \eta R(s, a, s') & \text{if } s' \text{ is terminal} \\ (1 - \eta)Q^*(s, a) + \eta [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')] & \text{otherwise} \end{cases}$$

- Better convergence in practice
- Watch out your memory usage!
- Space complexity for SARSA/ $Q$ -learning:  $O(|\mathcal{S}||\mathcal{A}|)$ 
  - Store  $Q^*(s, a)$ 's or  $Q_\pi(s, a)$ 's for all  $(s, a)$  combinations

# Remarks

- Update rule for  $Q^*$  (or  $Q_\pi$ ) takes terminal states differently:

$$Q^*(s, a) \leftarrow \begin{cases} (1 - \eta)Q^*(s, a) + \eta R(s, a, s') & \text{if } s' \text{ is terminal} \\ (1 - \eta)Q^*(s, a) + \eta [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')], & \text{otherwise} \end{cases}$$

- Better convergence in practice
- Watch out your memory usage!
- Space complexity for SARSA/ $Q$ -learning:  $O(|\mathbb{S}||\mathbb{A}|)$ 
  - Store  $Q^*(s, a)$ 's or  $Q_\pi(s, a)$ 's for all  $(s, a)$  combinations
- Why not train a (deep) regressor for  $Q^*(s, a)$ 's or  $Q_\pi(s, a)$ 's?
  - Space reduction due to the generalizability of the regressor