

# Deep Learning

## Lecture-16:

# Reinforcement Learning

*Datalab*

# Outline

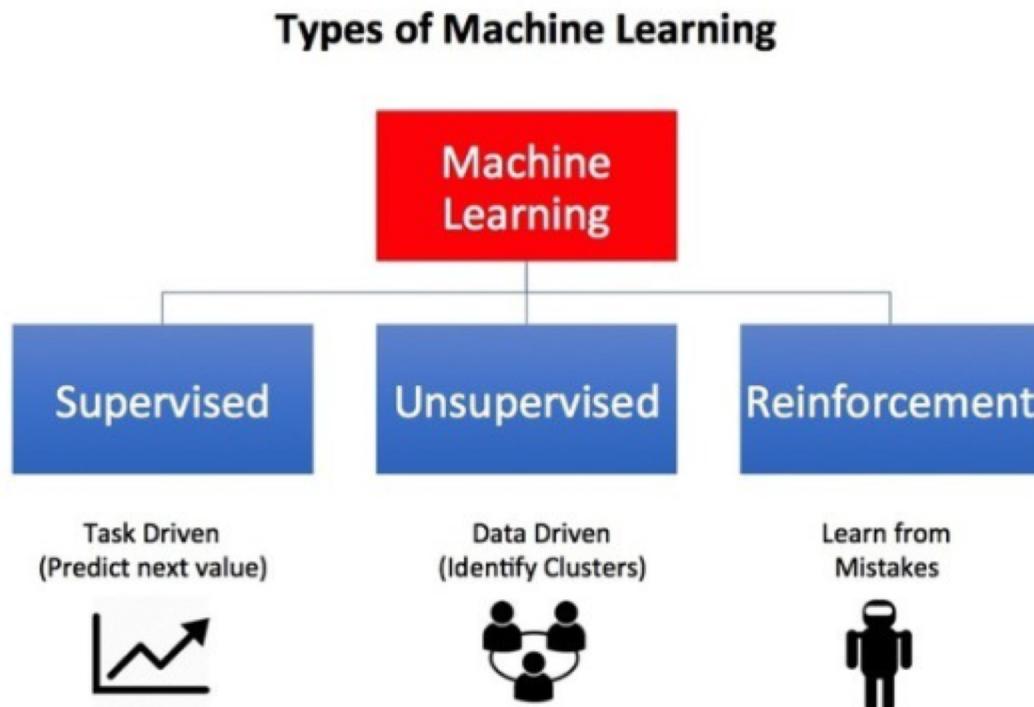
- Introduction of RL
- Markov Decision Process(MDP)
- Value Iteration
- Policy Iteration

# Outline

- Introduction of RL
- Markov Decision Process(MDP)
- Value Iteration
- Policy Iteration

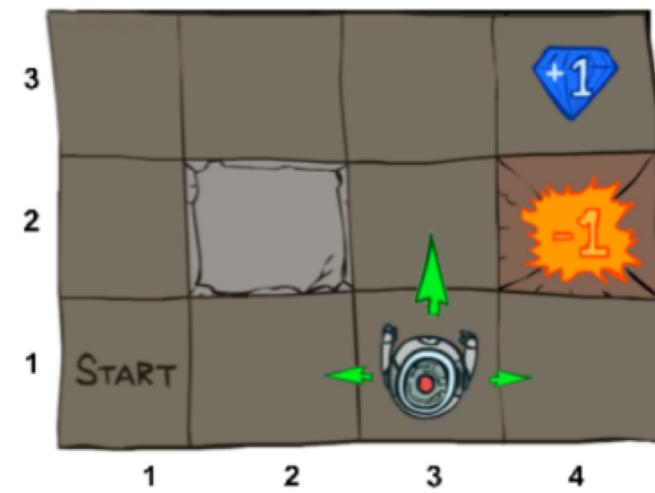
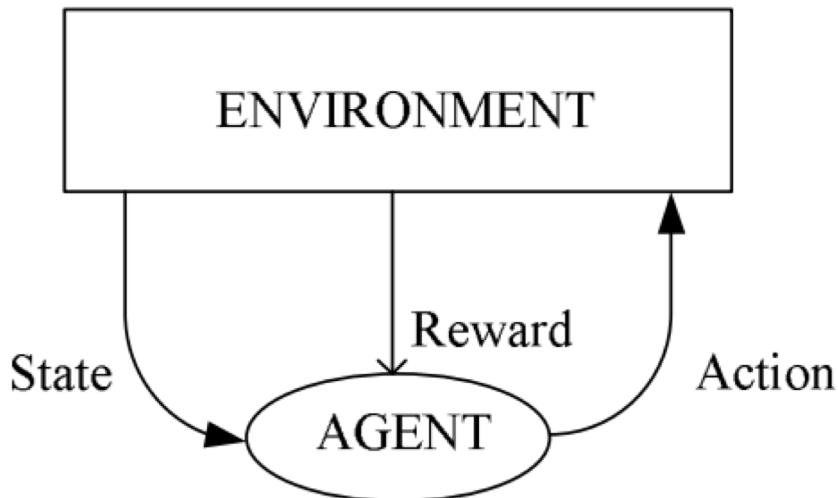
# Introduction of RL

- Types of Machine Learning



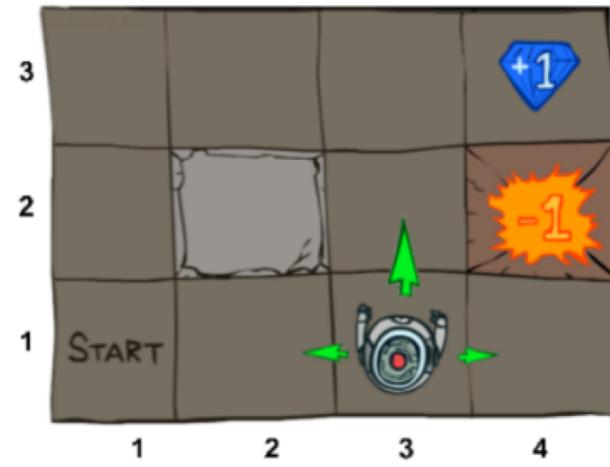
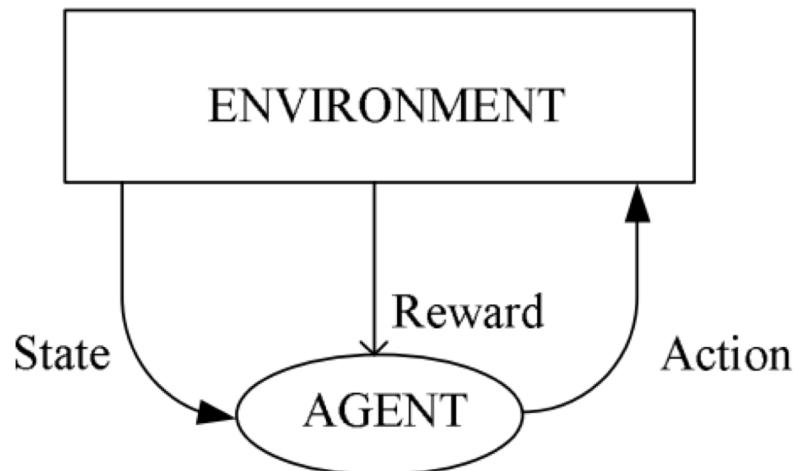
# Introduction of RL

- 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  $a^{(t)}, a^{(t-1)}, \dots$



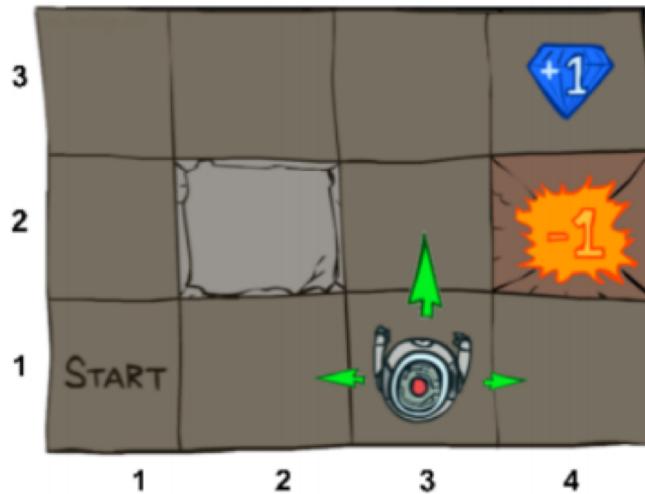
# Introduction of RL

- 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 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  $\mathbf{a}^{(t)}, \mathbf{a}^{(t-1)}, \dots$
- No **what** to predict ( $\mathbf{y}^{(i)}$ 's), just **how good** a prediction is ( $R^{(t)}$ 's)
  - $R^{(t)}$ 's are also called the **critics**



# Applications



Kohl and Stone, 2004



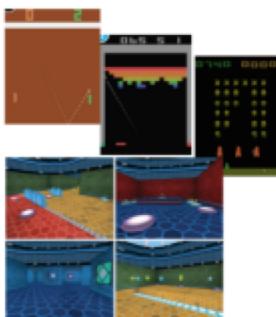
Ng et al, 2004



Tedrake et al, 2005

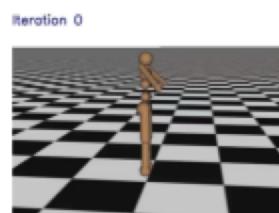


Kober and Peters, 2009



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)



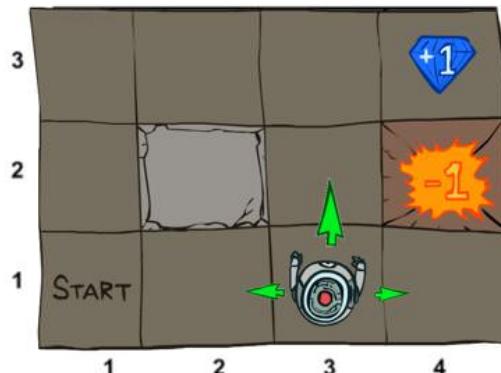
Silver\*, Huang\*, et  
al, 2016  
(AlphaGo)

# Outline

- Introduction of RL
- Markov Decision Process(MDP)
- Value Iteration
- Policy Iteration

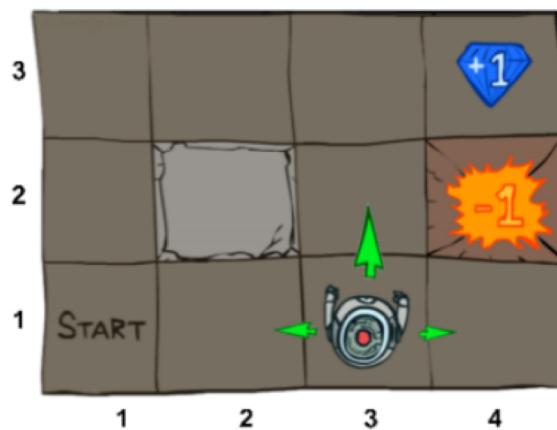
# MDP

- Definition:
  - 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” generally means that given the present state, the next state and the past are independent



# MDP

- 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



# MDP

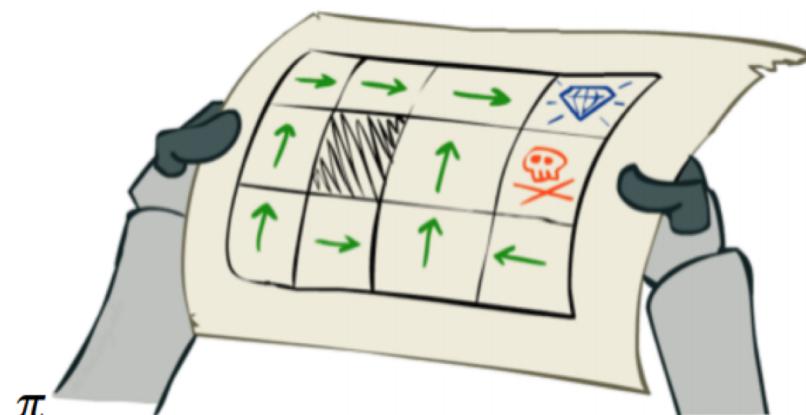
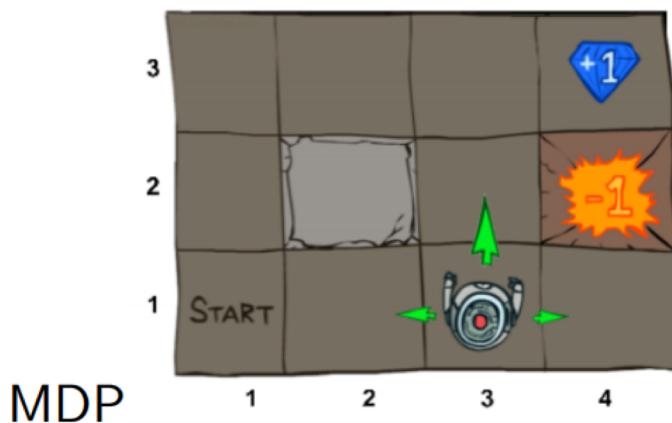
- Given a policy  $\pi(s) = \mathbf{a}$ , an MDP proceeds as follows:

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

with the accumulative reward

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

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



# Goal of MDP

- 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

- Introduction of RL
- Markov Decision Process(MDP)
- **Value Iteration**
- Policy Iteration

# Value iteration

- 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

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

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

# Value iteration

- Dynamic Programming

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

- $h = H - 1$ :

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

- $h = H - 2$ :

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

- $h = 0$ :

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

- $h = -1$ :

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

# Value iteration

- Dynamic Programming

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

- $h = H - 1$ :

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

- $h = H - 2$ :

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

- $h = 0$ :

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

- $h = -1$ :

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

 0 0 0	 -1 -1 -1	 -2 -2 -2	 -3 -2 -3
 0 0 0	 -1 -1 -1	 -2 -1 -2	 -2 -1 -2
 0  0	 -1 0 -1	 -1 0 -1	 -1 0 -1

初始化

第一轮迭代

第二轮迭代

第三轮迭代

# Value iteration

- Algorithm

**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**

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

**end**

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

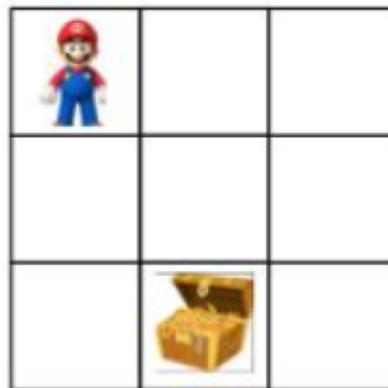
**foreach**  $s$  **do**

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

**end**

# Examples for value iteration

- Game: Super Mario
  - State: each grid(9 states)
  - Action: up, down, left, right(4 actions) in each state(excluding terminal state)
  - $P(s' | s, a) = 1$ (no noise)
  - Get -1 reward after each move
  - Terminal state(box): Reward = 0



# Examples for value iteration

- Game: Super Mario

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

	0	0
0	0	0
0		0

初始化

	-1	-1
-1	-1	-1
-1		-1

第一轮迭代

	-2	-2
-2	-1	-2
-1		-1

第二轮迭代

	-3	-3
-2	-1	-2
-1		-1

第三轮迭代

# Examples for value iteration

- Recall the algorithm

**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

0	0	0
0	0	0
0	0	0

初始化

1	-1	-1
-1	-1	-1
-1	0	-1

第一轮迭代

-2	-2	-2
-2	-1	-2
-1	0	-1

第二轮迭代

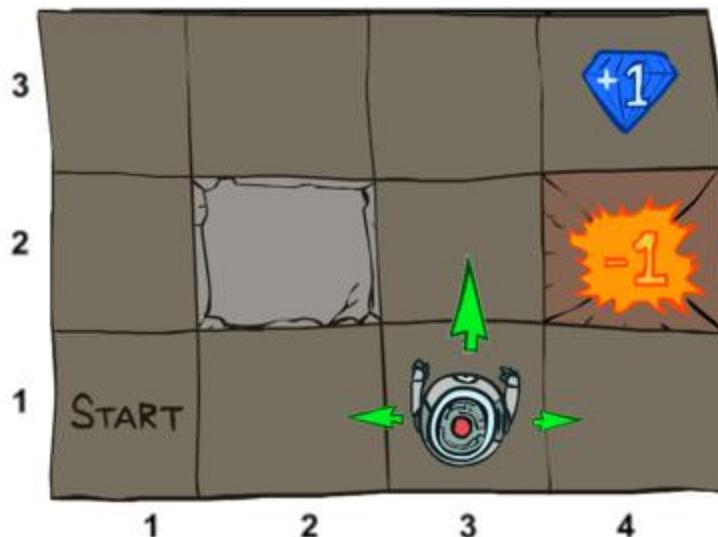
-3	-2	-3
-2	-1	-2
-1	0	-1

第三轮迭代

# Examples for value iteration

- Game: Gridworld

- Noise of transition probability  $P(s'|s; a)$ : 0.2
- $\gamma$ : 0.9

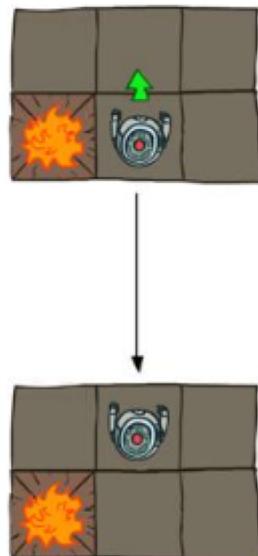


# Examples for value iteration

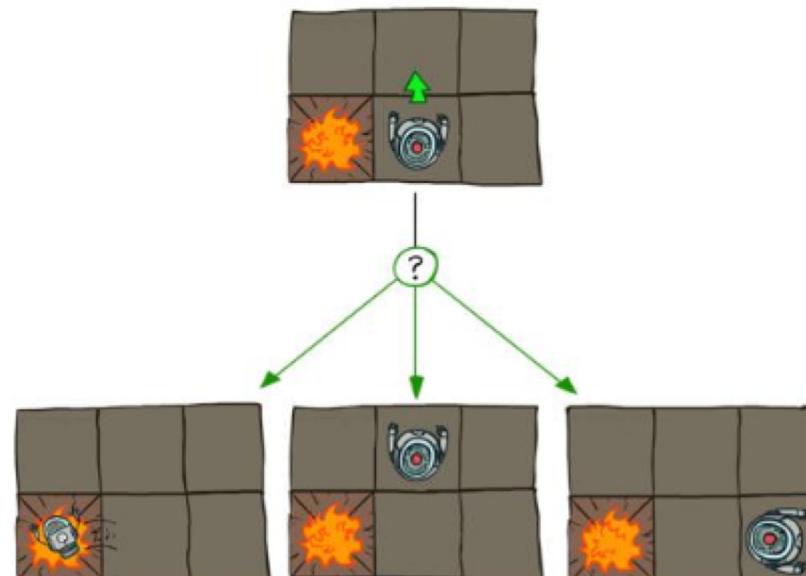
- Game: Gridworld
  - Noise of transition probability  $P(s'|s; a)$ : 0.2

## Grid World Actions

Deterministic Grid World

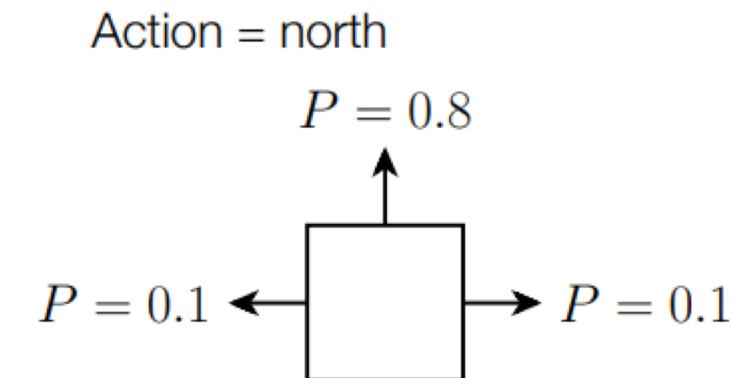
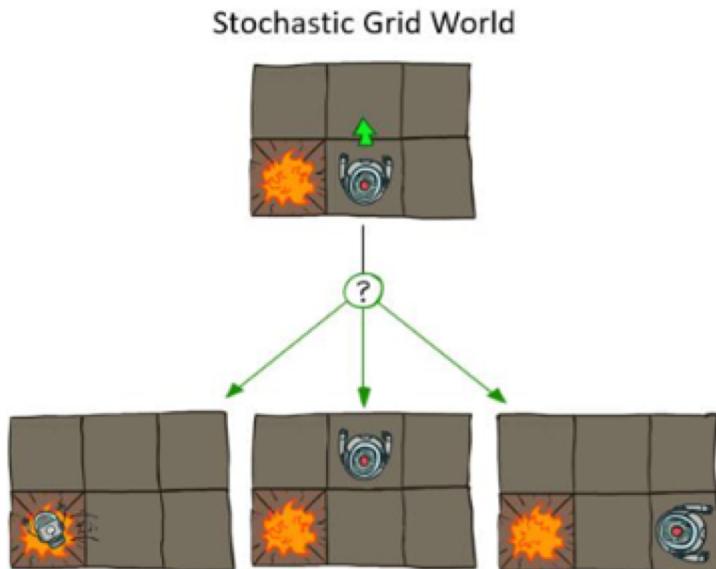


Stochastic Grid World



# Examples for value iteration

- Game: Gridworld
  - Noise: 0.2(0.1 go to the left, 0.1 go to the right)

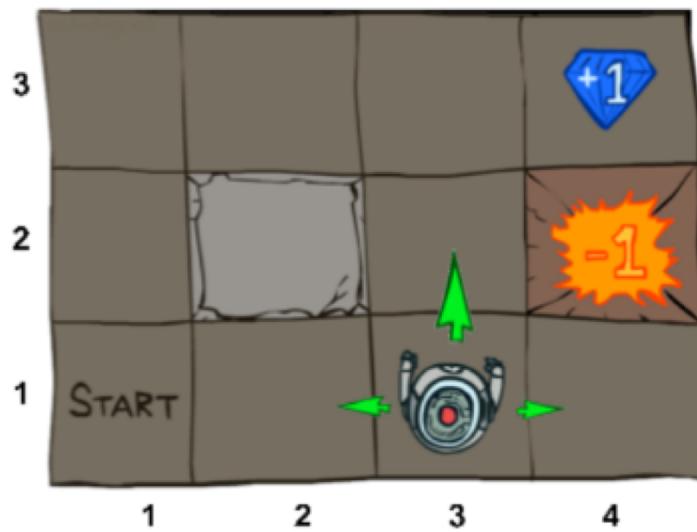


# Examples for value iteration

- Game: Gridworld

$$V_0(s) \leftarrow 0$$

$k = 0$



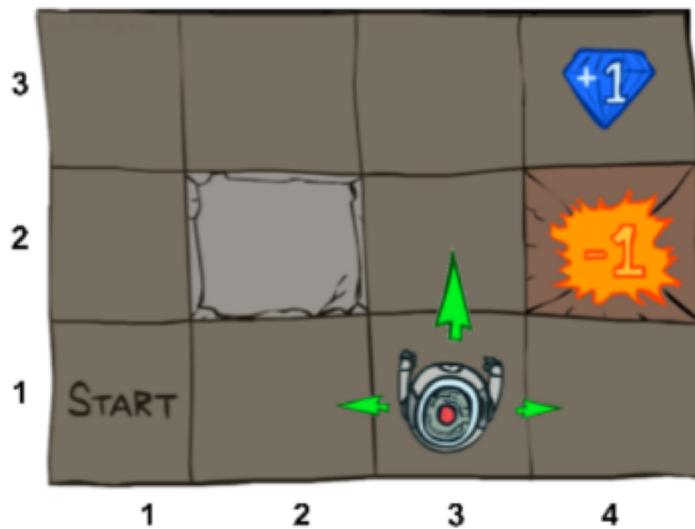
0.00	0.00	0.00	0.00
0.00		0.00	0.00
0.00	0.00	0.00	0.00

# Examples for value iteration

- Game: Gridworld

$$V_1(s) \leftarrow \max_a \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma V_0(s'))$$

$k = 1$



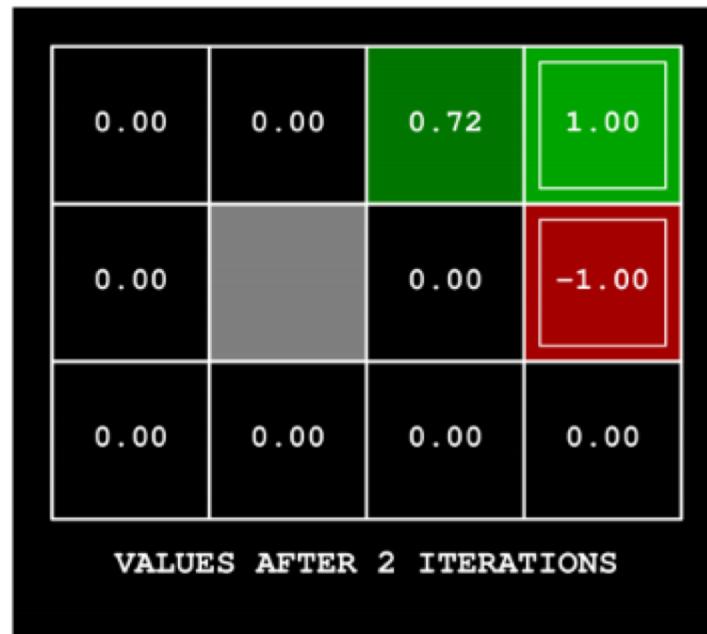
0.00	0.00	0.00	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

# Examples for value iteration

- Game: Gridworld

$$V_2(s) \leftarrow \max_a \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma V_1(s'))$$

$k = 2$



# Examples for value iteration

- Game: Gridworld

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma V_k(s'))$$

$k = 3$



# Examples for value iteration

- Game: Gridworld

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma V_k(s'))$$

$k = 4$



# Examples for value iteration

- Game: Gridworld

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma V_k(s'))$$

k = 100



# Examples for value iteration

- Game: Gridworld
  - Policy Extraction

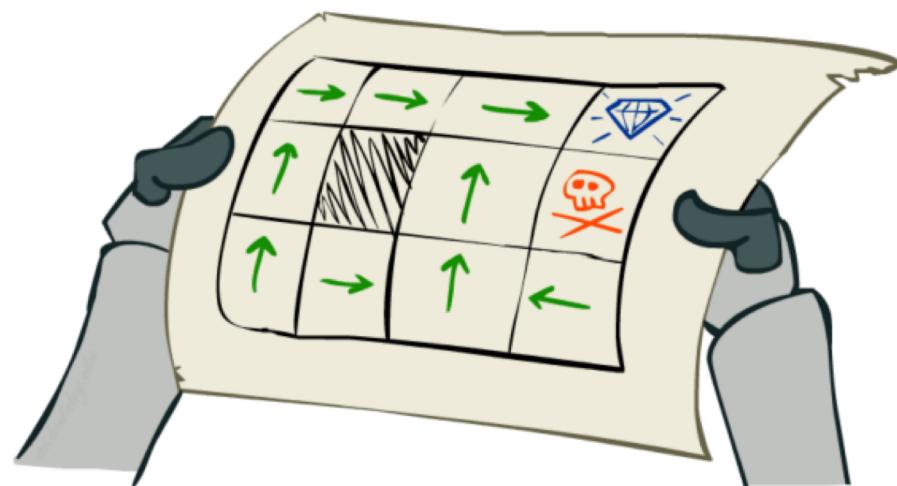
**foreach**  $s$  **do**

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

**end**

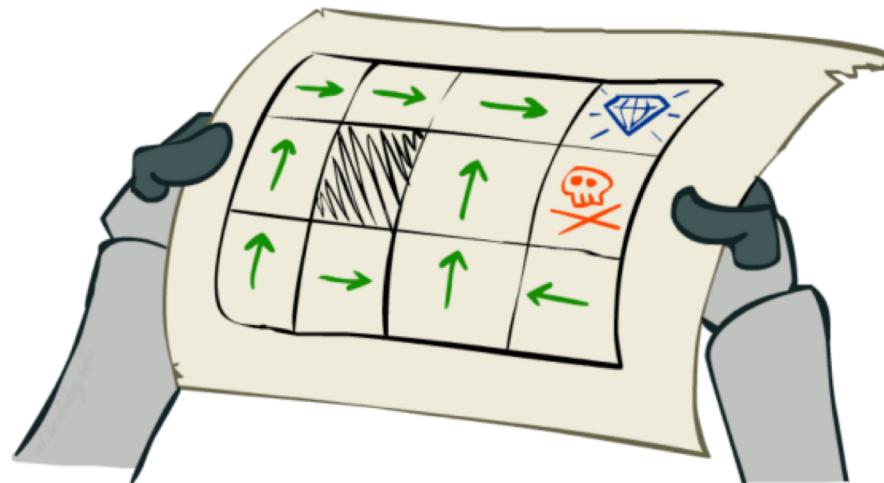
0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.43	0.48	0.28

VALUES AFTER 100 ITERATIONS



# Examples for value iteration

- Game: Gridworld
  - We get the optimal policy after 100 iterations
  - The policy  $\pi$  gives an action for each state
  - Following the optimal  $\pi$ , we can get the maximum expected total reward



# Examples for value iteration

- Theorem
  - Value iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$
  - Proof: Please see the slide of the teacher

# Outline

- Introduction of RL
- Markov Decision Process(MDP)
- Value Iteration
- **Policy Iteration**

# Policy iteration

- Goal
  - Given an MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$
  - 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

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

# Policy iteration

- Goal
  - Given an MDP  $(\mathbb{S}, \mathbb{A}, P, R, \gamma, H \rightarrow \infty)$
  - 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

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

# Policy iteration

- Policy Evaluation
- Policy Improvement

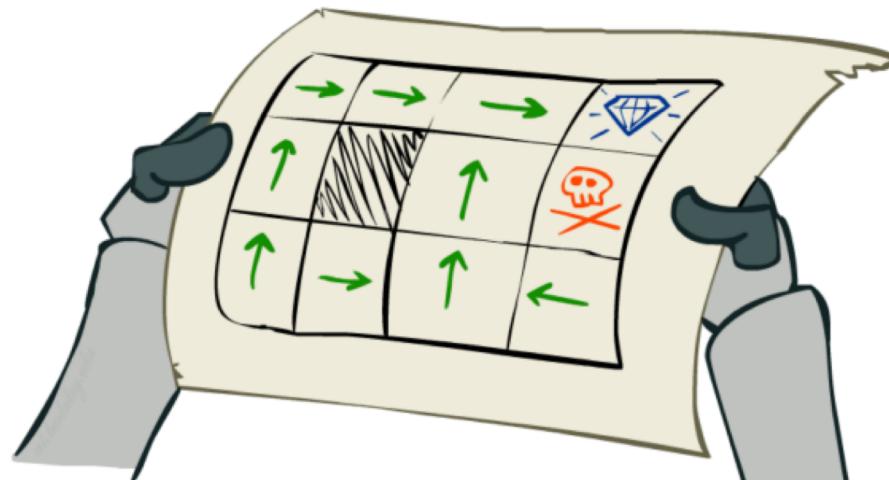
# Policy iteration

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

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

- Recall optimal value function in value iteration

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

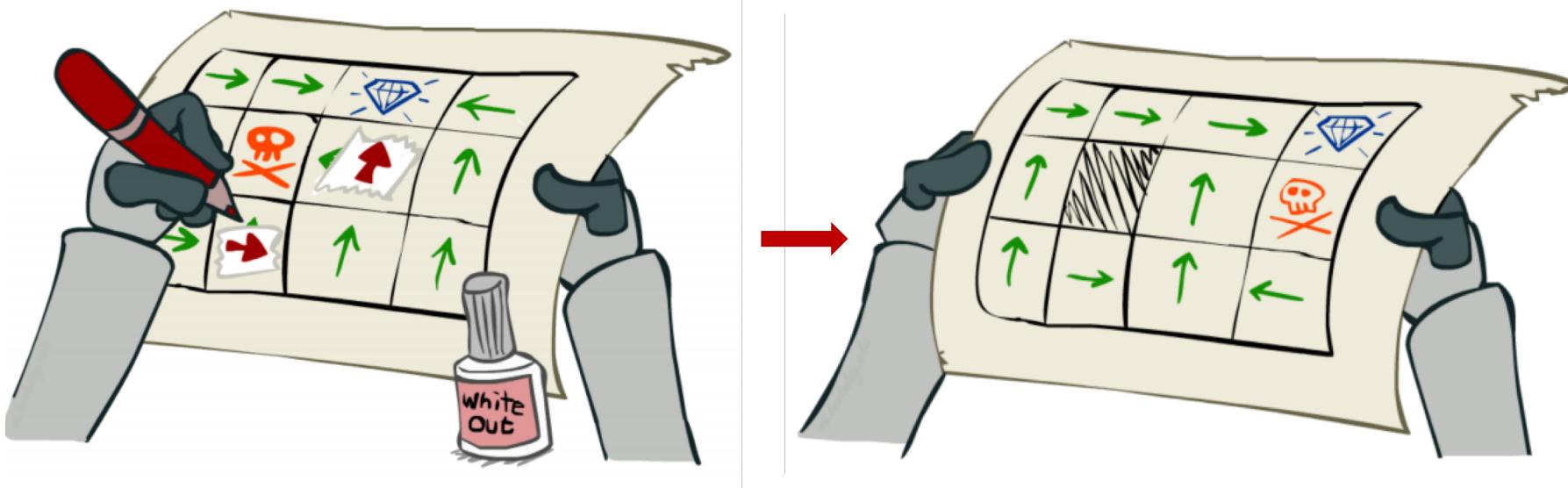


# Policy iteration

- Policy Improvement

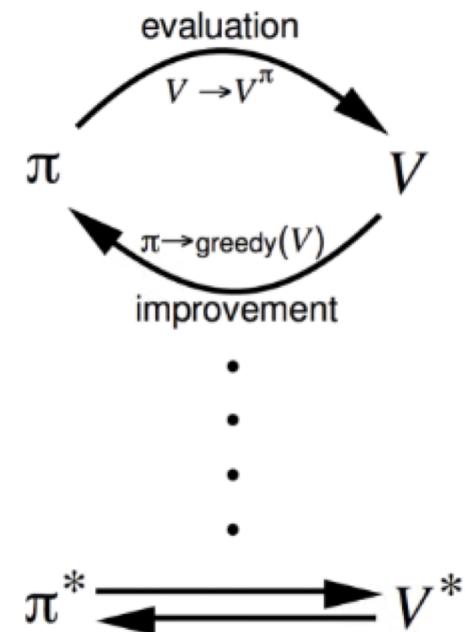
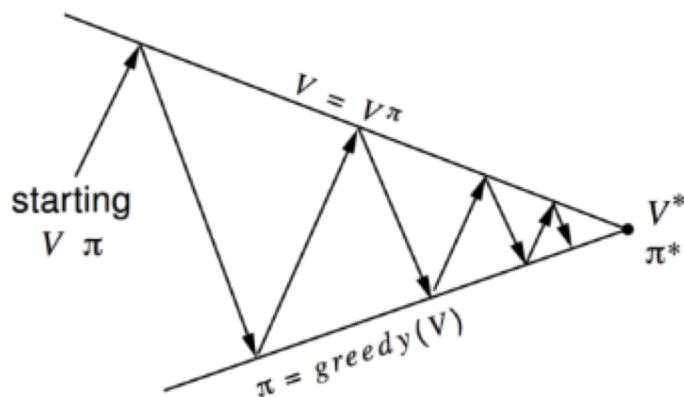
Update rule: for all  $s$  do

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



# Policy iteration

- Introduction



# Policy iteration

- Algorithm

**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**

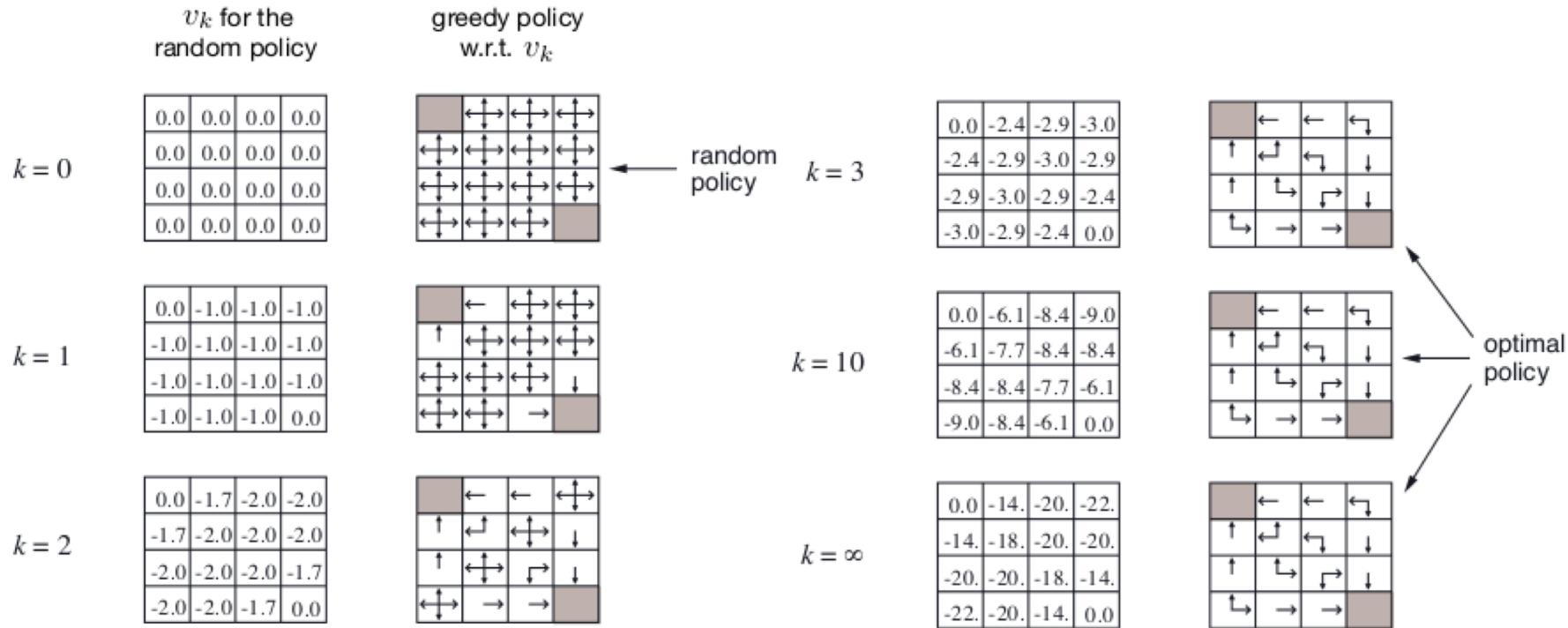
$| \pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')];$

    | **end**

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

# Examples for policy iteration

- Gridworld



# Examples for policy iteration

- Theorem
  - Policy iteration converges and gives the optimal policy  $\pi^*$  when  $H \rightarrow \infty$ .
  - Proof: Please see the slide of the teacher

# Value iteration VS. Policy iteration

- Difference:

```

1. Initialization
 $v(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$ 

2. Policy Evaluation
    Repeat
         $\Delta \leftarrow 0$ 
        For each  $s \in \mathcal{S}$ :
             $temp \leftarrow v(s)$ 
             $v(s) \leftarrow \sum_{s'} p(s'|s, \pi(s)) [r(s, \pi(s), s') + \gamma v(s')]$ 
             $\Delta \leftarrow \max(\Delta, |temp - v(s)|)$ 
        until  $\Delta < \theta$  (a small positive number)

3. Policy Improvement
     $policy-stable \leftarrow true$ 
    For each  $s \in \mathcal{S}$ :
         $temp \leftarrow \pi(s)$ 
         $\pi(s) \leftarrow \arg \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$ 
        If  $temp \neq \pi(s)$ , then  $policy-stable \leftarrow false$ 
    If  $policy-stable$ , then stop and return  $v$  and  $\pi$ ; else go to 2

```

**finding optimal value function**

Initialize array  $v$  arbitrarily (e.g.,  $v(s) = 0$  for all  $s \in \mathcal{S}^+$ )

Repeat

$\Delta \leftarrow 0$   
 For each  $s \in \mathcal{S}$ :  
 $temp \leftarrow v(s)$   
 $v(s) \leftarrow \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$   
 $\Delta \leftarrow \max(\Delta, |temp - v(s)|)$   
 until  $\Delta < \theta$  (a small positive number)

Output a deterministic policy,  $\pi$ , such that

$\pi(s) = \arg \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$

Figure 4.5: Value iteration.

**one policy update (extract policy from the optimal value function)**

Figure 4.3: Policy iteration (using iterative policy evaluation) for  $v_*$ . This algorithm has a subtle bug, in that it may never terminate if the policy continually switches between two or more policies that are equally good. The bug can be fixed by adding additional flags, but it makes the pseudocode so ugly that it is not worth it. :-)

# Value iteration VS. Policy iteration

- Similarity:
  - Both are Model-based methods
  - We need to know  $P(s' | s, a)$  and  $R(s, a, s')$

# Next class

- We need to know  $P(s'|s, a)$  and  $R(s, a, s')$  to use value iteration and policy iteration
- In the real world, the transition probability and reward function may be unknown, what should we do to get the optimal policy?

# Thanks!