

# Reinforcement Learning I

Aarti Singh

Machine Learning 10-701

Apr 3, 2023

Slides courtesy: Henry Chai, Eric Xing



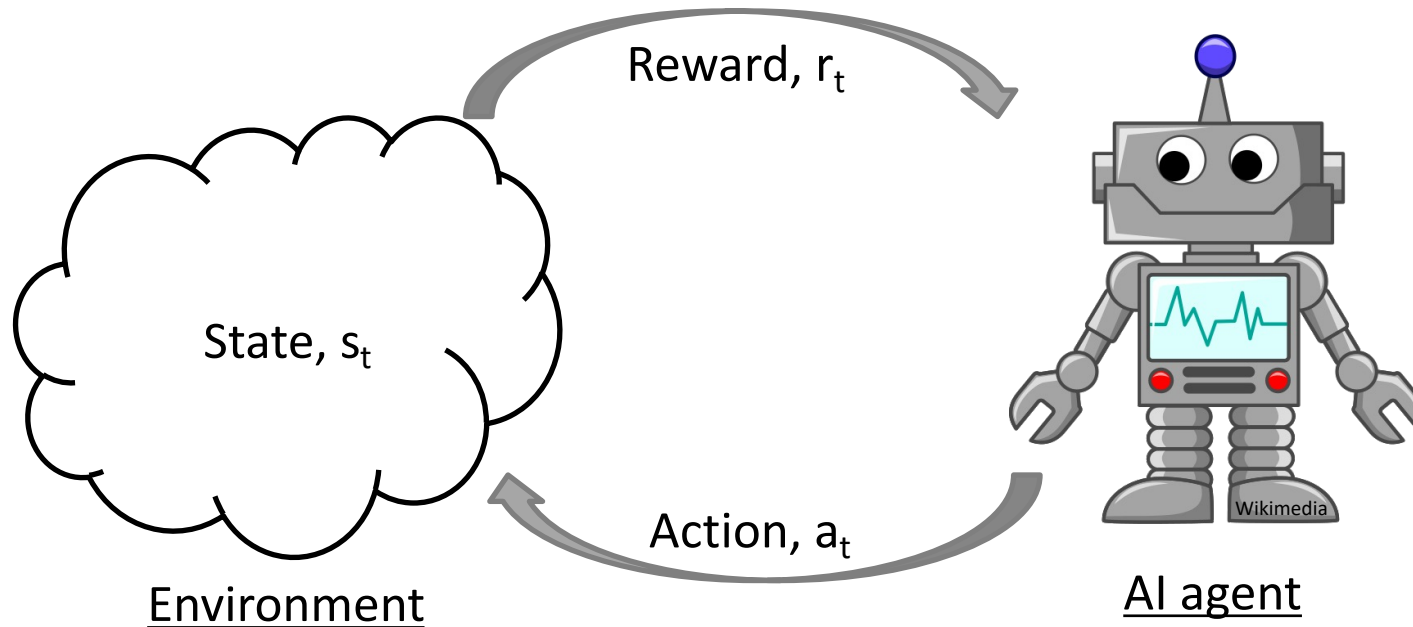
**MACHINE LEARNING** DEPARTMENT

**Carnegie Mellon.**  
School of Computer Science

# Learning Tasks

- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ 
  - Regression -  $y^{(i)} \in \mathbb{R}$
  - Classification -  $y^{(i)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ 
  - Clustering
  - Dimensionality reduction
- Reinforcement learning -  $\mathcal{D} = \{\mathbf{s}^{(t)}, \mathbf{a}^{(t)}, r^{(t)}\}_{t=1}^T$

# RL setup



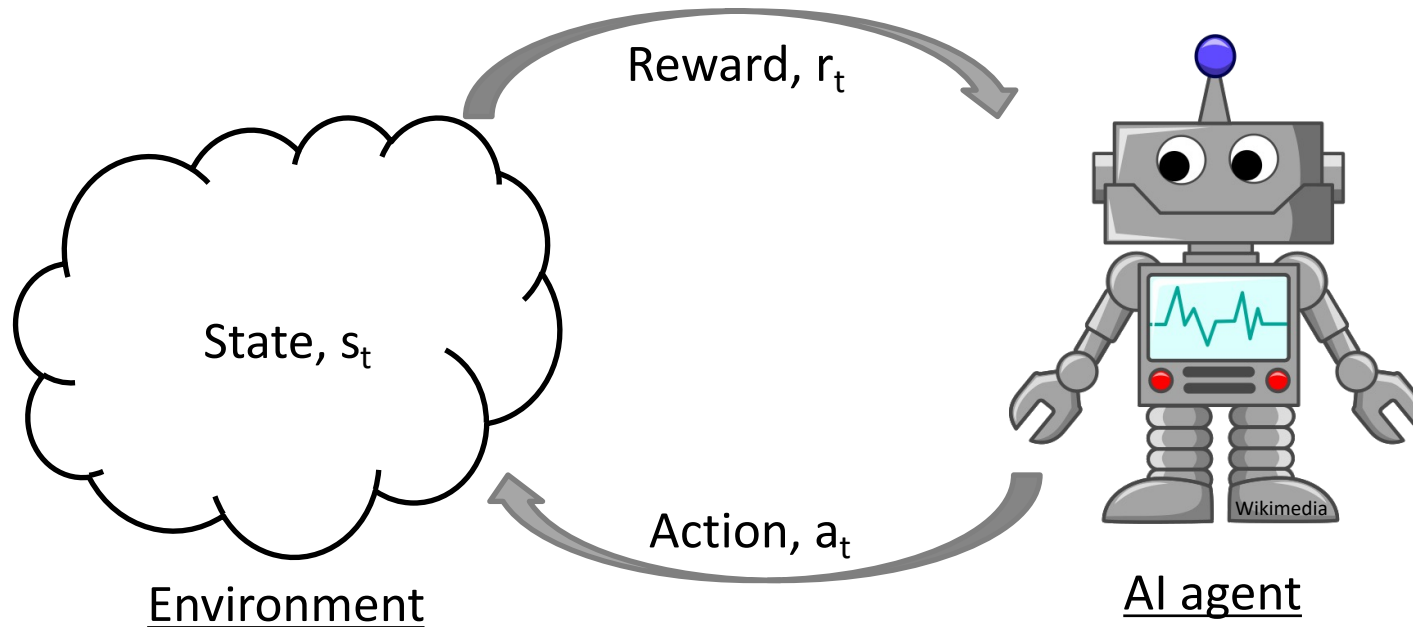
Agent chooses **actions** which can depend on past

Environment can change **state** with each action

**Reward** (Output) depends on (Inputs) action and state of environment

Goal: Maximize total reward

# Differences from supervised learning



- Maximize reward (rather than learn reward)
- Inputs are not iid – state & action depends on past
- Can control some inputs - actions

# RL examples



<https://techobserver.net/2019/06/argo-ai-self-driving-car-research-center/>



<https://www.cmu.edu/news/stories/archives/2017/september/snakebot-mexico.html>



<https://www.wired.com/2012/02/high-speed-trading/>



<https://twitter.com/alphgomovie>

# RL setup

- State space,  $\mathcal{S}$
- Action space,  $\mathcal{A}$
- Reward function
  - Stochastic,  $p(r | s, a)$
  - Deterministic,  $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Transition function
  - Stochastic,  $p(s' | s, a)$
  - Deterministic,  $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$
- Reward and transition functions can be known or unknown

# RL setup

- Policy,  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ 
  - Specifies an action to take in *every* state
- Value function,  $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ 
  - Measures the expected total reward of starting in some state  $s$  and executing policy  $\pi$ , i.e., in every state, taking the action that  $\pi$  returns

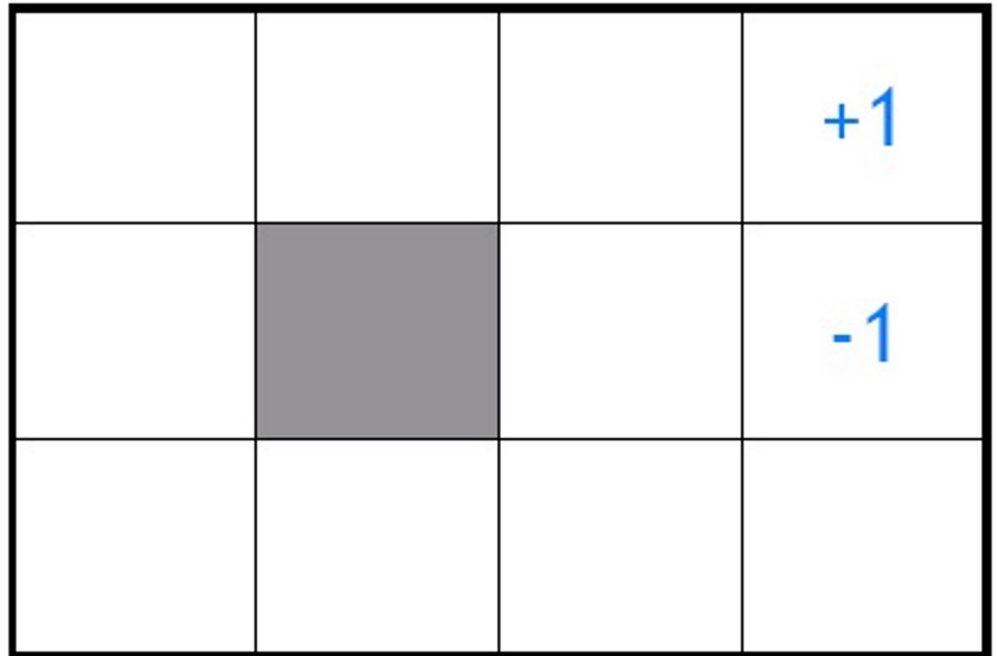
# RL example - gridworld

$\mathcal{S}$  = all empty squares in the grid

$\mathcal{A}$  = {up, down, left, right}

Deterministic transitions

Rewards of +1 and -1 for entering the labelled squares



Terminate after receiving either reward



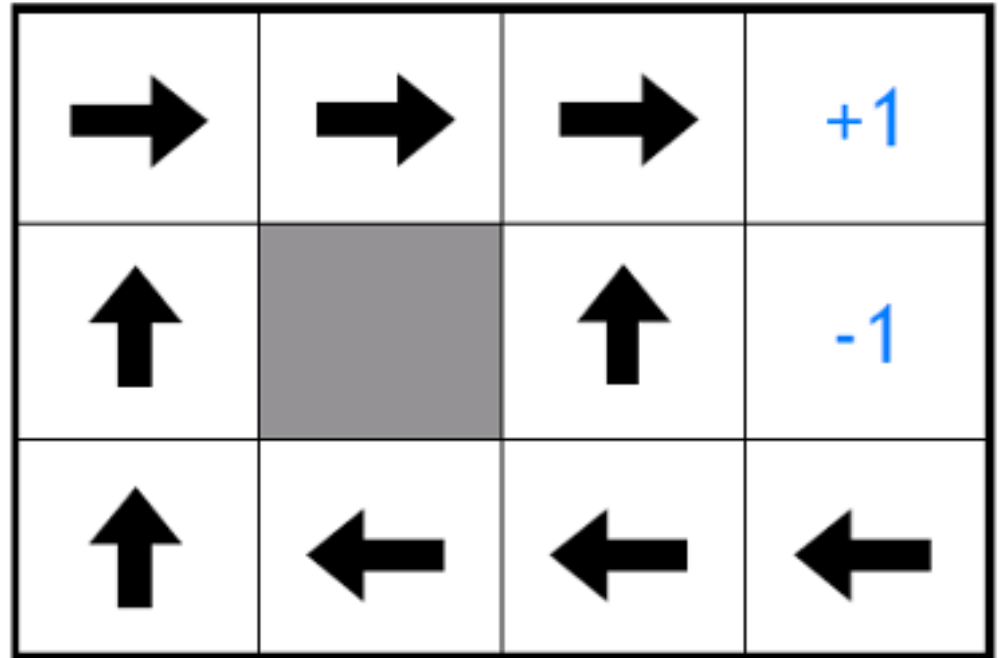
# RL example - gridworld

$\mathcal{S}$  = all empty squares in the grid

$\mathcal{A}$  = {up, down, left, right}

Deterministic transitions

Rewards of +1 and -1 for entering the labelled squares



Poll: Is this policy optimal?

Terminate after receiving either reward

# RL example - gridworld

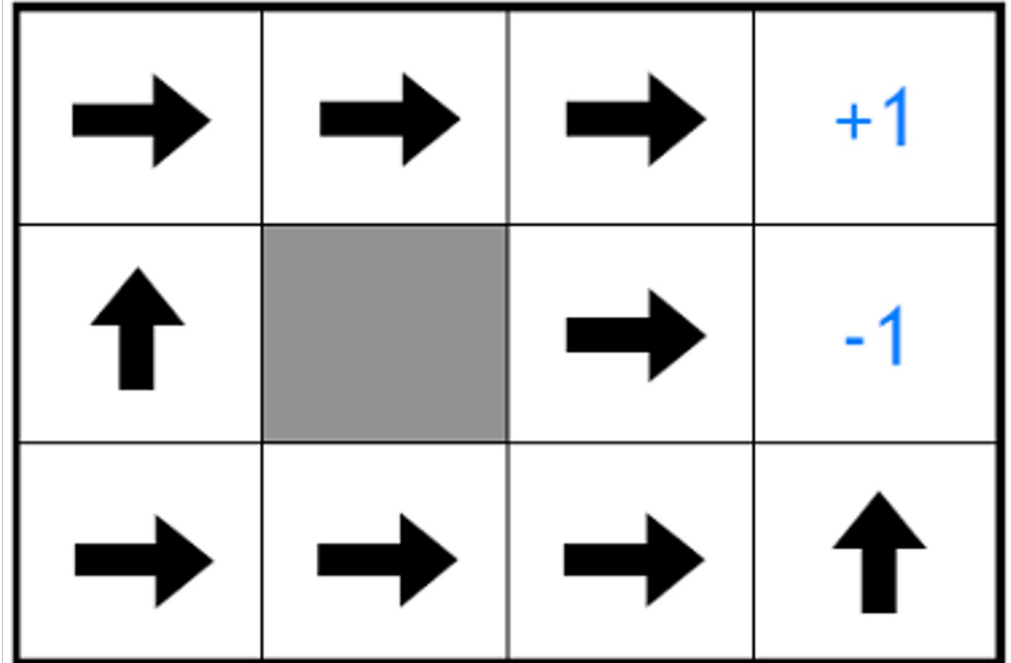
$\mathcal{S}$  = all empty squares in the grid

$\mathcal{A}$  = {up, down, left, right}

Deterministic transitions

Rewards of +1 and -1 for entering the labelled squares

Terminate after receiving either reward



Optimal policy given a reward of  
-2 per step

# RL example - gridworld

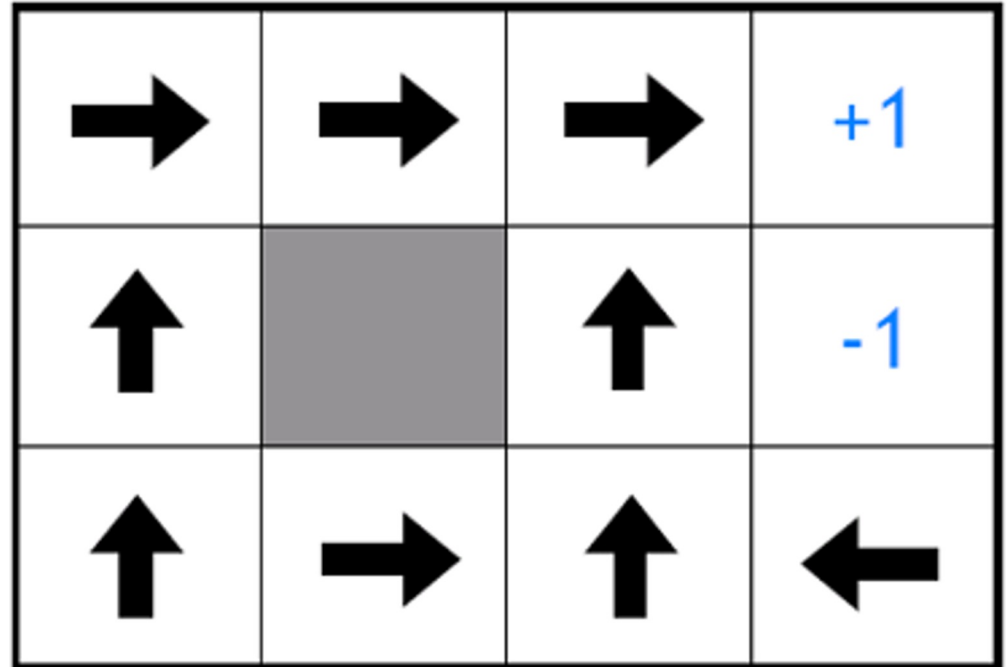
$\mathcal{S}$  = all empty squares in the grid

$\mathcal{A}$  = {up, down, left, right}

Deterministic transitions

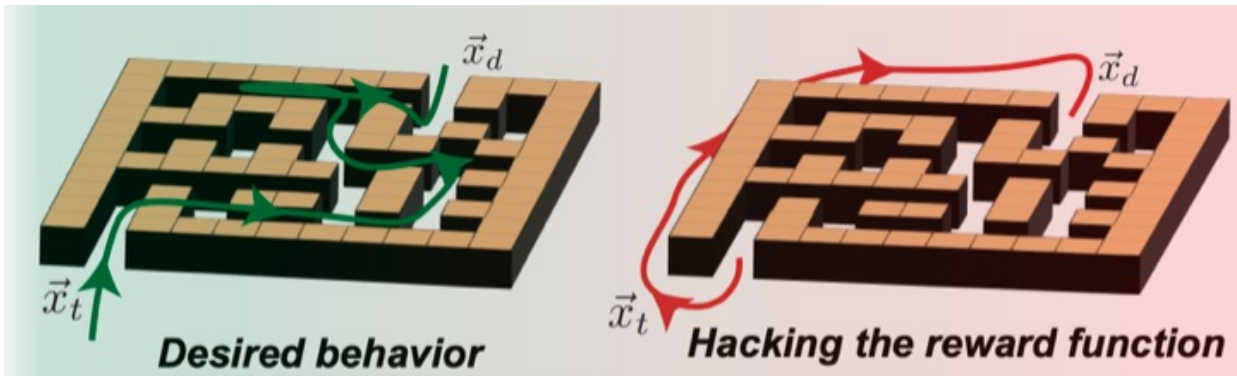
Rewards of +1 and -1 for entering the labelled squares

Terminate after receiving either reward



Optimal policy given a reward of  
-0.1 per step

# Reward hacking



Alhub.org

[Amodei-Clark'16]



# Markov Decision Process

1. Start in some initial state  $s_0$
  2. For time step  $t$ :
    - a. Agent observes state  $s_t$
    - b. Agent takes action  $a_t = \pi(s_t)$       Deterministic policy
    - c. Agent receives reward  $r_t \sim p(r | s_t, a_t)$
    - d. Agent transitions to state  $s_{t+1} \sim p(s' | s_t, a_t)$
- MDPs make the **Markov assumption**: the reward and next state only depend on the current state and action.

# Discounted Reward

Total reward is  $\sum_{t=0}^{\infty} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \dots$

where  $0 < \gamma < 1$  is some discount factor for future rewards

Why discount?

- Mathematically tractable – total reward doesn't explode

$$1 + 1 + 1 + \dots = \infty \text{ but } 1 + 0.8*1 + (0.8)^2*1 + \dots = 5$$

- Risk aversion under uncertainty
- Actions don't have lasting impact

# Key challenges

- The algorithm has to gather its own training data
- The outcome of taking some action is often stochastic or unknown until after the fact
- Decisions can have a delayed effect on future outcomes (exploration-exploitation tradeoff)
  - explore** decisions whose reward is uncertain
  - exploit** decisions which give high reward

# MDP example: Multi-armed bandits

Single state:  $|\mathcal{S}| = 1$

Three actions:  $\mathcal{A} = \{1, 2, 3\}$

Deterministic transitions

Rewards are stochastic





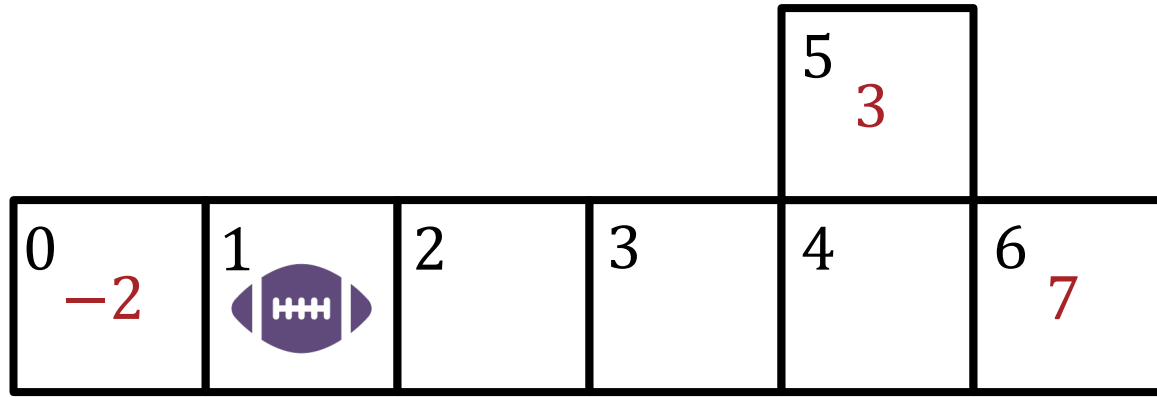


# RL: objective function

- Find a policy  $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(s) \quad \forall s \in \mathcal{S}$
- $V^{\pi}(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E} [R(s_0 = s, \pi(s_0))$   
 $\quad + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots]$   
 $= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[R(s_t, \pi(s_t))]$

where  $0 < \gamma < 1$  is some discount factor for future rewards

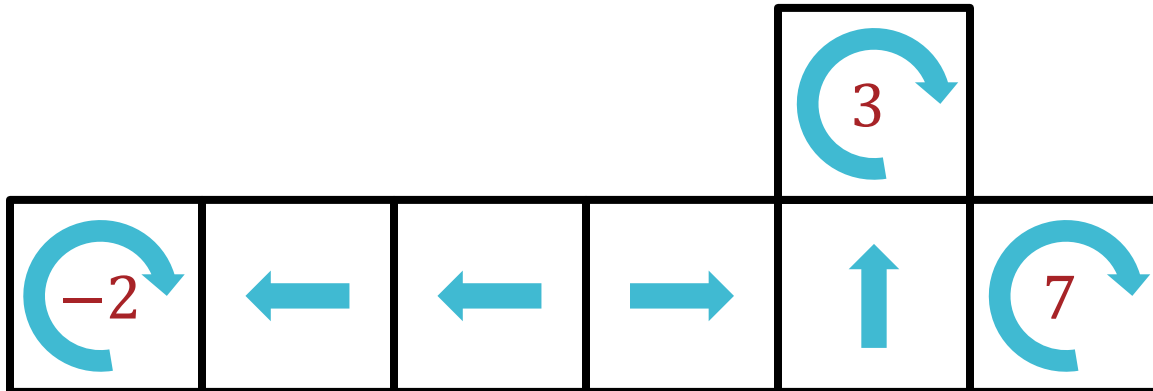
# Value function: example



$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

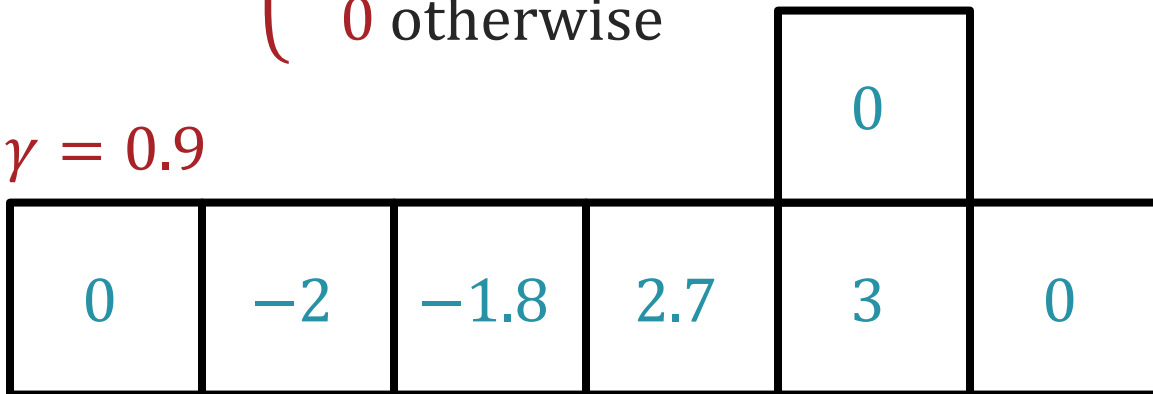
$$\gamma = 0.9$$

# Value function: example

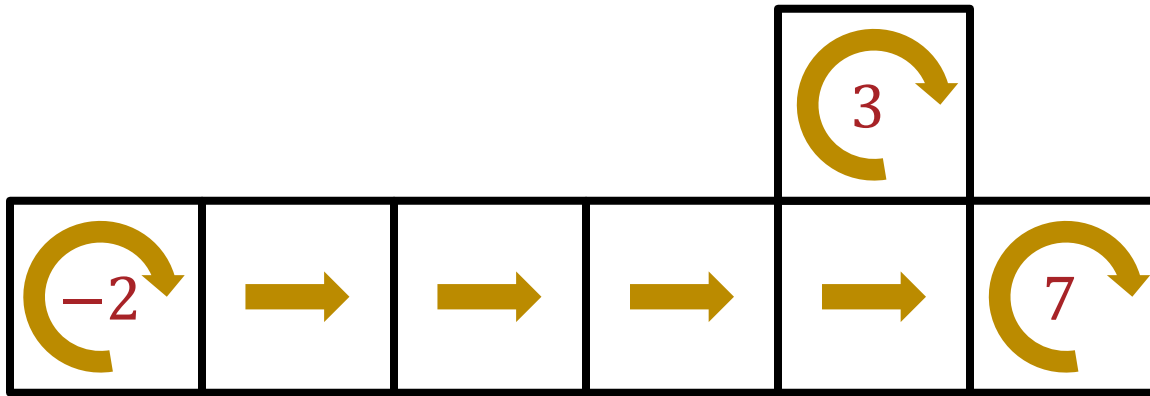


$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

$$\gamma = 0.9$$

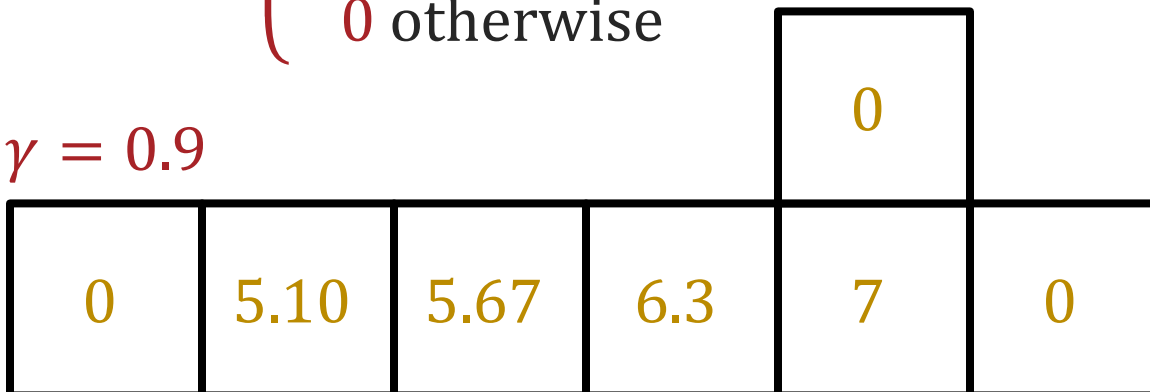


# Value function: example



$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

$$\gamma = 0.9$$



# Value function – deterministic reward

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) V^\pi(s_1)$$

**Bellman equations**

# Optimal value function and policy

- Optimal value function:

$$V^*(s) = \max_{a \in \mathcal{A}} [R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')]$$

- System of  $|\mathcal{S}|$  equations and  $|\mathcal{S}|$  variables – nonlinear!

- Optimal policy:

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} \underbrace{R(s, a)}_{\text{Immediate reward}} + \underbrace{\gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')}_{\text{Expected (Discounted) Future reward}}$$

- Insight: if you know the optimal value function, you can solve for the optimal policy!

# Value iteration

- Inputs:  $R(s, a)$ ,  $p(s' | s, a)$ ,  $0 < \gamma < 1$
- Initialize  $V^{(0)}(s) = 0 \forall s \in \mathcal{S}$  (or randomly) and set  $t = 0$
- While not converged, do:

- For  $s \in \mathcal{S}$

$$V^{(t+1)}(s) \leftarrow \max_{a \in \mathcal{A}} \underbrace{[R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{(t)}(s')]}_{Q(s, a)}$$

- $t = t + 1$

- For  $s \in \mathcal{S}$

$$\pi^*(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} [R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{(t)}(s')]$$

- Return  $\pi^*$



# Value iteration

- Inputs:  $R(s, a)$ ,  $p(s' | s, a)$ ,  $0 < \gamma < 1$
- Initialize  $V^{(0)}(s) = 0 \forall s \in \mathcal{S}$  (or randomly) and set  $t = 0$
- While not converged, do:

- For  $s \in \mathcal{S}$ 
  - For  $a \in \mathcal{A}$

$$Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^{(t)}(s')$$

- $V^{(t+1)}(s) \leftarrow \max_{a \in \mathcal{A}} Q(s, a)$

- $t = t + 1$

- For  $s \in \mathcal{S}$ 
  - $\pi^*(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$

- Return  $\pi^*$

# Poll

- What is the runtime per iteration?
  - A.  $O(1)$
  - B.  $|S| |A|$
  - C.  $|S| |A|^2$
  - D.  $|A| |S|^2$
  - E.  $|A|^2 |S|^2$

# Value iteration: convergence

- Runtime per iteration:  $O(|S|^2|A|)$

**Theorem 1:** Value function convergence

$V$  will converge to  $V^*$  if each state is “visited” infinitely often (Bertsekas, 1989)

**Theorem 2:** Convergence criterion

if  $\max_{s \in \mathcal{S}} |V^{(t+1)}(s) - V^{(t)}(s)| < \epsilon$ ,  
then  $\max_{s \in \mathcal{S}} |V^{(t+1)}(s) - V^*(s)| < \frac{2\epsilon\gamma}{1-\gamma}$  (Williams & Baird, 1993)

**Theorem 3:** Policy convergence

The “greedy” policy,  $\pi(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$ , converges to the optimal  $\pi^*$  in a finite number of iterations, often before the value function has converged! (Bertsekas, 1987)

# Policy iteration

➤ Can we learn the policy directly, instead of first learning the value function?

- Inputs:  $R(s, a)$ ,  $p(s' | s, a)$ ,  $0 < \gamma < 1$

- Initialize  $\pi$  randomly

- While not converged, do:

- Solve the Bellman equations defined by policy  $\pi$

Now linear!

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) V^\pi(s')$$

- Update  $\pi$

$$\pi(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^\pi(s')$$

- Return  $\pi$

# Policy iteration: convergence

- Runtime per iteration:  $O(|S|^2|A| + |S|^3)$

# Poll

- How many policies are there?
  - A.  $|S| + |A|$
  - B.  $|S| |A|$
  - C.  $|S|^{|A|}$
  - D.  $|A|^{|S|}$

# Policy iteration: convergence

- Runtime per iteration:  $O(|S|^2|A| + |S|^3)$
- Number of policies:  $|S| |A|$
- Policy improves each iteration
- Thus, the number of iterations needed to converge is bounded!
- Empirically, policy iteration requires fewer iterations than value iteration.

# Next Questions

- How to handle unknown state transition and reward functions?
- How to handle continuous states and actions?