

# Reinforcement Learning II

Aarti Singh

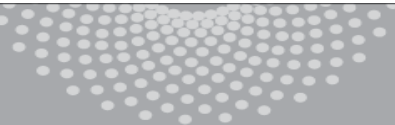
Machine Learning 10-701

Apr 5, 2023

Slides courtesy: Henry Chai, Eric Xing

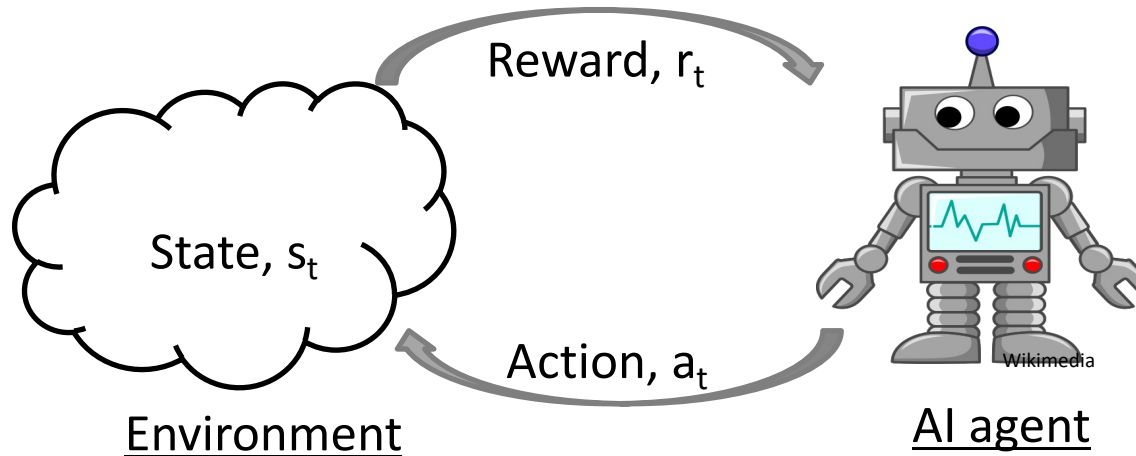


**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# RL setup



1. Start in some initial state  $s_0$
2. For time step  $t$ :
  - a. Agent observes state  $s_t$
  - b. Agent takes action  $a_t = \pi(s_t)$
  - c. Agent receives reward  $r_t \sim p(r \mid s_t, a_t)$
  - d. Agent transitions to state  $s_{t+1} \sim p(s' \mid s_t, a_t)$

# RL setup

- Policy,  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ 
  - Specifies an action to take in *every* state
- Value function,  $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ 
  - $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[R(s_t, \pi(s_t))]$$
 R – deterministic reward
- **Goal:** Find policy that maximizes expected discounted total reward

$$\pi^* = \operatorname{argmax}_{\pi} V^\pi(s) \quad \forall s \in \mathcal{S}$$

# Bellman Equation

Value function satisfies the set of recursive equations:

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

- 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:

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

# 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  $s \in \mathcal{S}$

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

- $t = t + 1$

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

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

- 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  $\pi^*$

# Value iteration: convergence

**Theorem 1:** Value function convergence

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

**Theorem 2:** Convergence criterion

$$\text{if } \max_{s \in \mathcal{S}} |V^{(t+1)}(s) - V^{(t)}(s)| < \epsilon,$$

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

**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

- Number of policies:  $|A|^{|S|}$
- 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?

# Optimal Q function and policy

- Deterministic rewards
- $Q^*(s, a) = \mathbb{E}[\text{total discounted reward of taking action } a \text{ in state } s, \text{ assuming all future actions are optimal}]$

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

$$V^*(s') = \max_{a' \in \mathcal{A}} Q^*(s', a')$$

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

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$$

- Insight: if we know  $Q^*$ , we can compute an optimal policy  $\pi^*$ !

# Optimal Q function and policy

- Deterministic rewards and state transitions
- $Q^*(s, a) = \mathbb{E}[\text{total discounted reward of taking action } a \text{ in state } s, \text{ assuming all future actions are optimal}]$

$$= R(s, a) + \gamma V^*(\delta(s, a))$$

- $V^*(\delta(s, a)) = \max_{a' \in \mathcal{A}} Q^*(\delta(s, a), a')$

$$Q^*(s, a) = R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q^*(\delta(s, a), a')$$

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$$

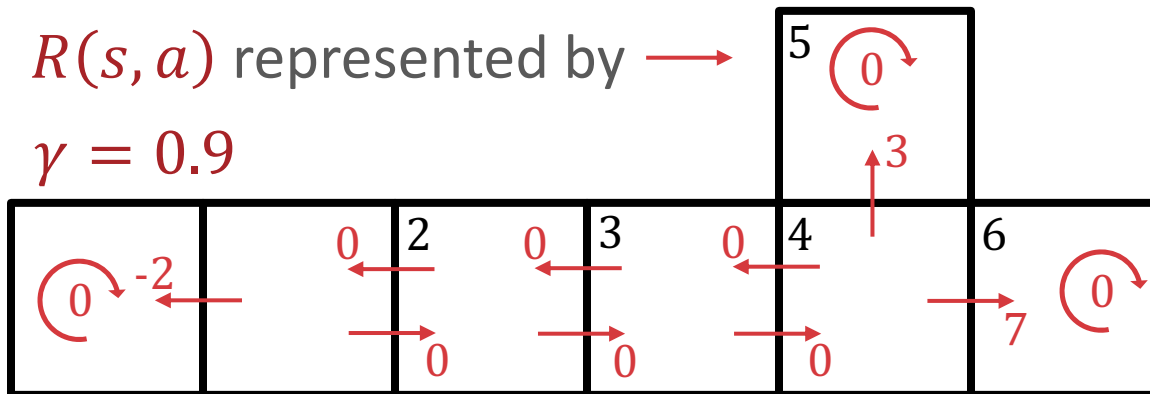
- Insight: if we know  $Q^*$ , we can compute an optimal policy  $\pi^*$ !

# Online Q-learning

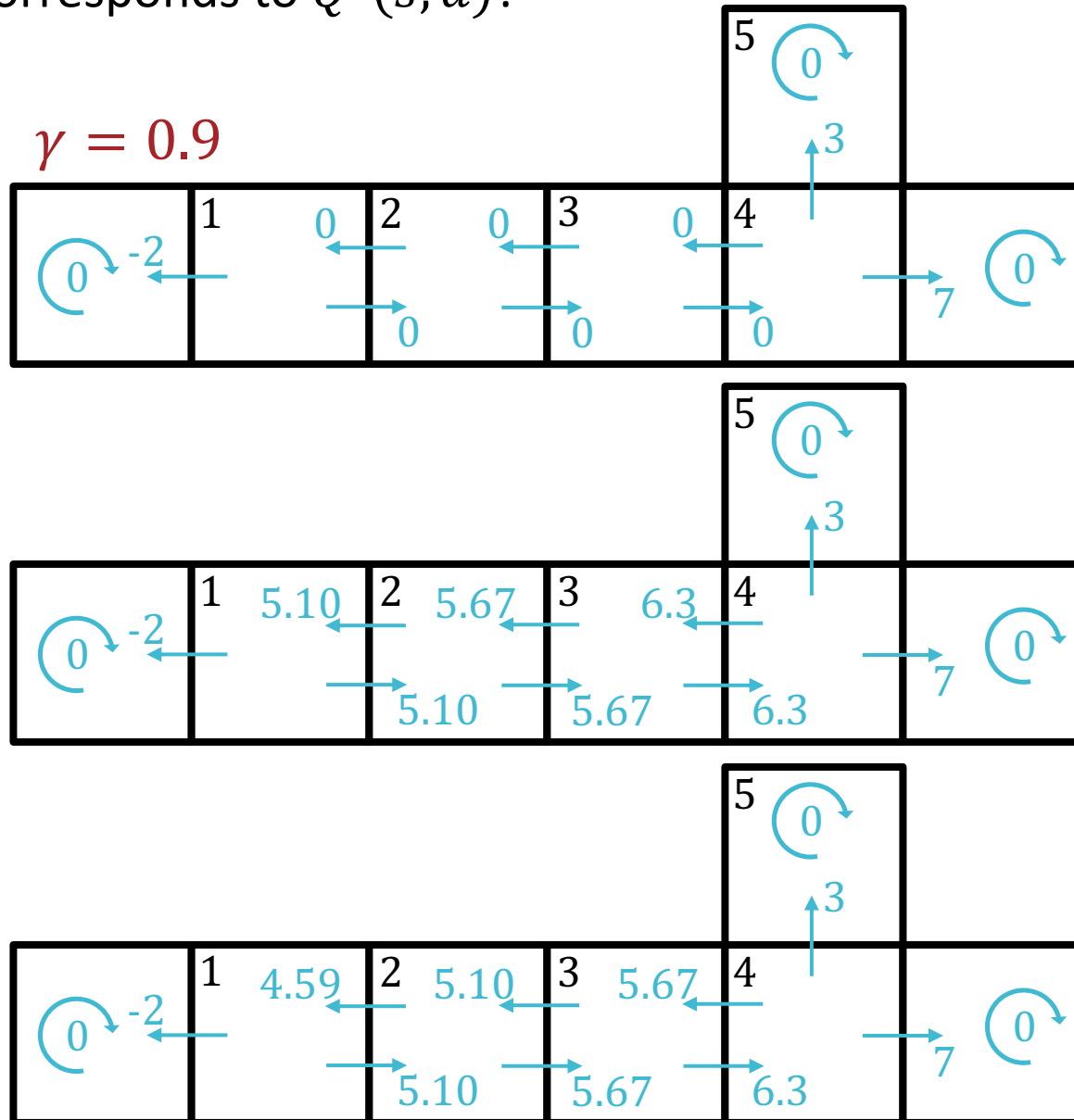
- Inputs: discount factor  $\gamma$ , an initial state  $s$
- Initialize  $Q(s, a) = 0 \forall s \in \mathcal{S}, a \in \mathcal{A}$  ( $Q$  is a  $|\mathcal{S}| \times |\mathcal{A}|$  array)
- While TRUE, do
  - Take a random action  $a$
  - Receive reward  $r = R(s, a)$
  - Update the state:  $s \leftarrow s'$  where  $s' = \delta(s, a)$
  - Update  $Q(s, a)$ :

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$

# Q-learning example



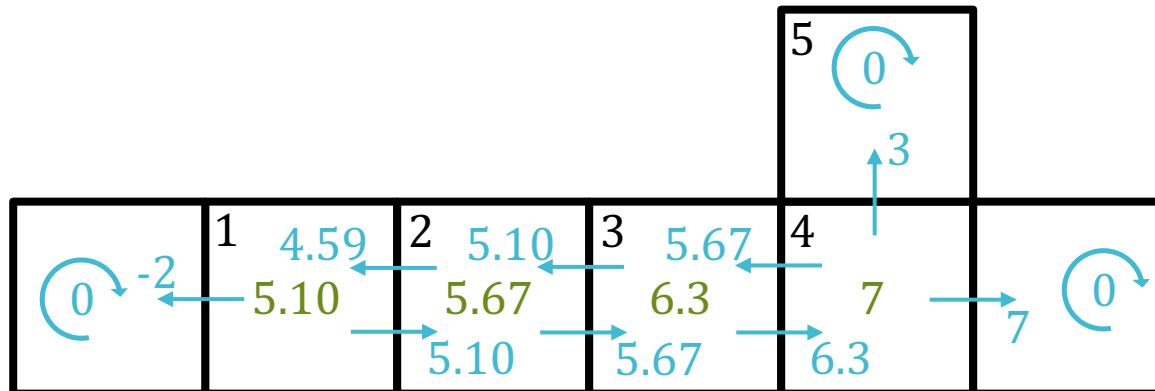
Which set of blue arrows  
 (roughly) corresponds to  $Q^*(s, a)$ ?



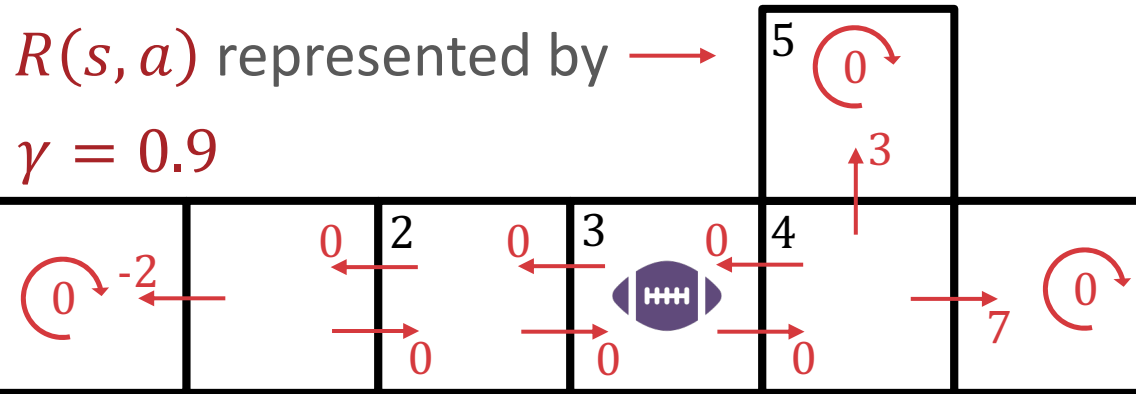
Which set of blue arrows  
(roughly) corresponds to  $Q^*(s, a)$ ?

$$Q^*(s, a) = R(s, a) + \gamma V^*(\delta(s, a))$$

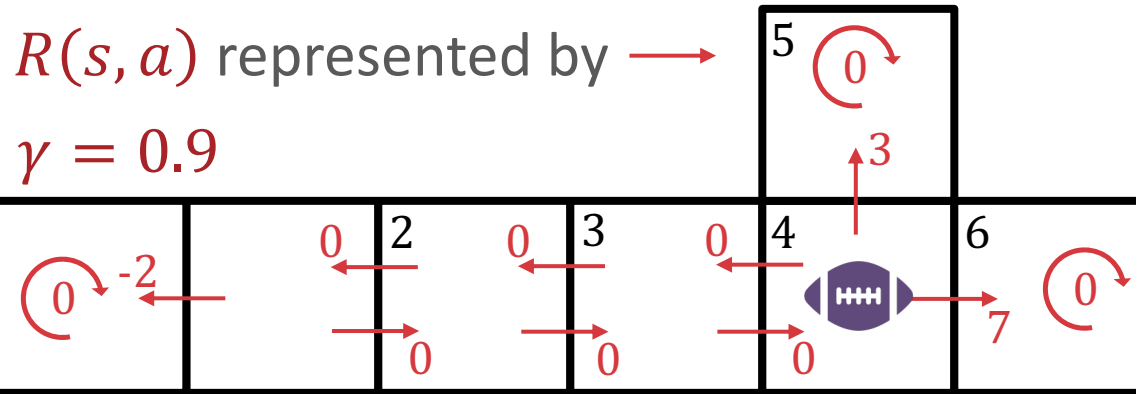
$V^*(s)$  shown in green





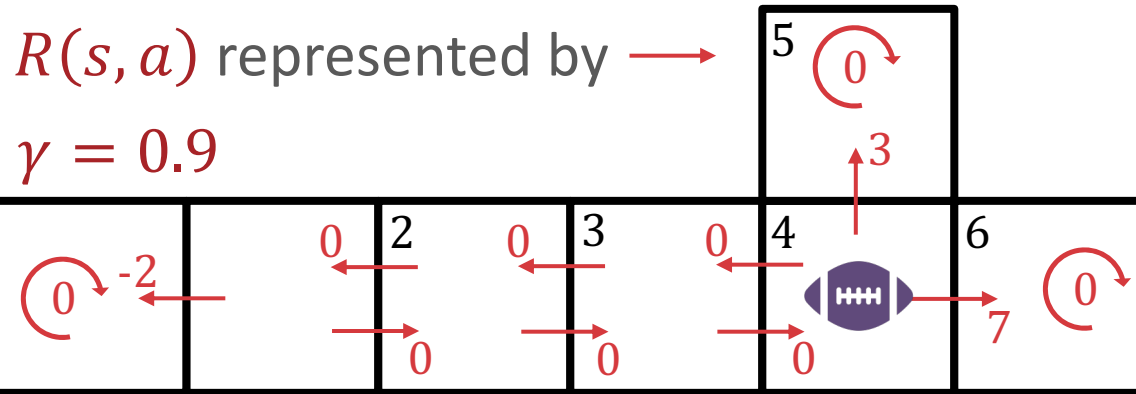


$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\curvearrowright$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0

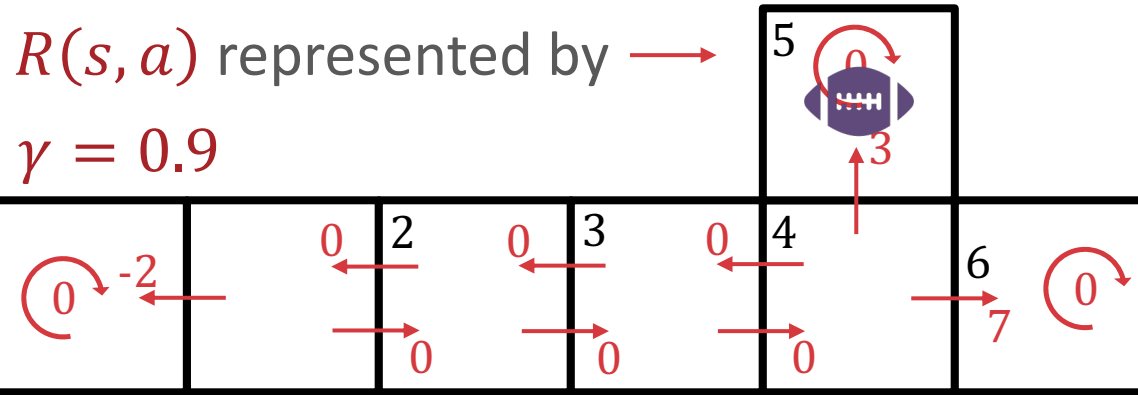


$$Q(3, \rightarrow) \leftarrow 0 + (0.9) \max_{a' \in \{\rightarrow, \leftarrow, \uparrow, \cup\}} Q(4, a') = 0$$

$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\cup$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0

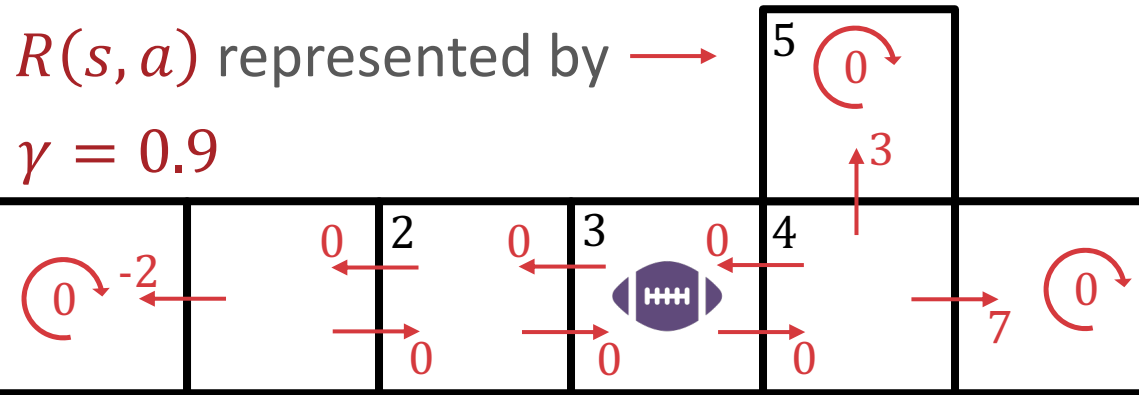


$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\curvearrowright$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0



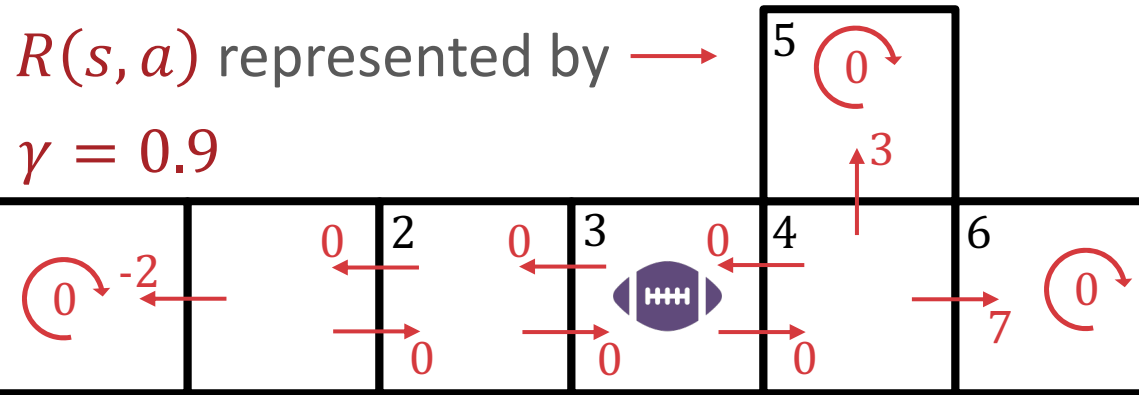
$$Q(4, \uparrow) \leftarrow 3 + (0.9) \max_{a' \in \{\rightarrow, \leftarrow, \uparrow, \cup\}} Q(5, a') = 3$$

$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\cup$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0



$$Q(3, \rightarrow) \leftarrow 0 + (0.9) \max_{a' \in \{\rightarrow, \leftarrow, \uparrow, \cup\}} Q(4, a') = 2.7$$

$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\cup$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	3	0
5	0	0	0	0
6	0	0	0	0



$$Q(3, \rightarrow) \leftarrow 0 + (0.9) \max_{a' \in \{\rightarrow, \leftarrow, \uparrow, \cup\}} Q(4, a') = 2.7$$

$Q(s, a)$	$\rightarrow$	$\leftarrow$	$\uparrow$	$\cup$
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	2.7	0	0	0
4	0	0	3	0
5	0	0	0	0
6	0	0	0	0

# Online Q-learning

- Inputs: discount factor  $\gamma$ , an initial state  $s$
- Initialize  $Q(s, a) = 0 \forall s \in \mathcal{S}, a \in \mathcal{A}$  ( $Q$  is a  $|\mathcal{S}| \times |\mathcal{A}|$  array)
- While TRUE, do
  - Take a random action  $a$
  - Receive reward  $r = R(s, a)$
  - Update the state:  $s \leftarrow s'$  where  $s' = \delta(s, a)$
  - Update  $Q(s, a)$ :

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$

# $\epsilon$ -greedy Online Q-learning

- Inputs: discount factor  $\gamma$ , an initial state  $s$ , greediness parameter  $\epsilon \in [0, 1]$
- Initialize  $Q(s, a) = 0 \forall s \in \mathcal{S}, a \in \mathcal{A}$  ( $Q$  is a  $|\mathcal{S}| \times |\mathcal{A}|$  array)
- While TRUE, do
  - With probability  $\epsilon$ , take the greedy action
$$a = \operatorname{argmax}_{a' \in \mathcal{A}} Q(s, a')$$
  - Otherwise, with probability  $1 - \epsilon$ , take a random action  $a$
  - Receive reward  $r = R(s, a)$
  - Update the state:  $s \leftarrow s'$  where  $s' = \delta(s, a)$
  - Update  $Q(s, a)$ :

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



# Stochastic Transitions

- Inputs: discount factor  $\gamma$ , an initial state  $s$ , greediness parameter  $\epsilon \in [0, 1]$ , learning rate  $\alpha \in [0, 1]$  (“trust parameter”)
- Initialize  $Q(s, a) = 0 \forall s \in \mathcal{S}, a \in \mathcal{A}$  ( $Q$  is a  $|\mathcal{S}| \times |\mathcal{A}|$  array)
- While TRUE, do
  - With probability  $\epsilon$ , take the greedy action
$$a = \operatorname{argmax}_{a' \in \mathcal{A}} Q(s, a')$$
Otherwise, with probability  $1 - \epsilon$ , take a random action  $a$
  - Receive reward  $r = R(s, a)$
  - Update the state:  $s \leftarrow s'$  where  $s' \sim p(s' | s, a)$
  - Update  $Q(s, a)$ :

$$Q(s, a) \leftarrow (1 - \alpha) \underbrace{Q(s, a)}_{\text{Current value}} + \alpha \underbrace{\left( r + \gamma \max_{a'} Q(s', a') \right)}_{\text{Update w/ deterministic transitions}}$$

# Temporal Difference Learning

- Inputs: discount factor  $\gamma$ , an initial state  $s$ , greediness parameter  $\epsilon \in [0, 1]$ , learning rate  $\alpha \in [0, 1]$  (“trust parameter”)
- Initialize  $Q(s, a) = 0 \forall s \in \mathcal{S}, a \in \mathcal{A}$  ( $Q$  is a  $|\mathcal{S}| \times |\mathcal{A}|$  array)
- While TRUE, do

- With probability  $\epsilon$ , take the greedy action

$$a = \operatorname{argmax}_{a' \in \mathcal{A}} Q(s, a')$$

Otherwise, with probability  $1 - \epsilon$ , take a random action  $a$

- Receive reward  $r = R(s, a)$
- Update the state:  $s \leftarrow s'$  where  $s' \sim p(s' | s, a)$
- Update  $Q(s, a)$ :

$$Q(s, a) \leftarrow \underbrace{Q(s, a)}_{\text{Current value}} + \alpha \left( \underbrace{r + \gamma \max_{a'} Q(s', a')}_{\text{Temporal difference target}} - Q(s, a) \right)$$

Temporal difference

# Q – learning: convergence

- For Algorithms 1 & 2 (deterministic transitions),  $Q$  converges to  $Q^*$  if
  1. Every valid state-action pair is visited infinitely often
    - Q-learning is exploration-insensitive: any visitation strategy that satisfies this property will work!
  2.  $0 \leq \gamma < 1$
  3.  $\exists \beta$  s.t.  $|R(s, a)| < \beta \forall s \in \mathcal{S}, a \in \mathcal{A}$
  4. Initial  $Q$  values are finite

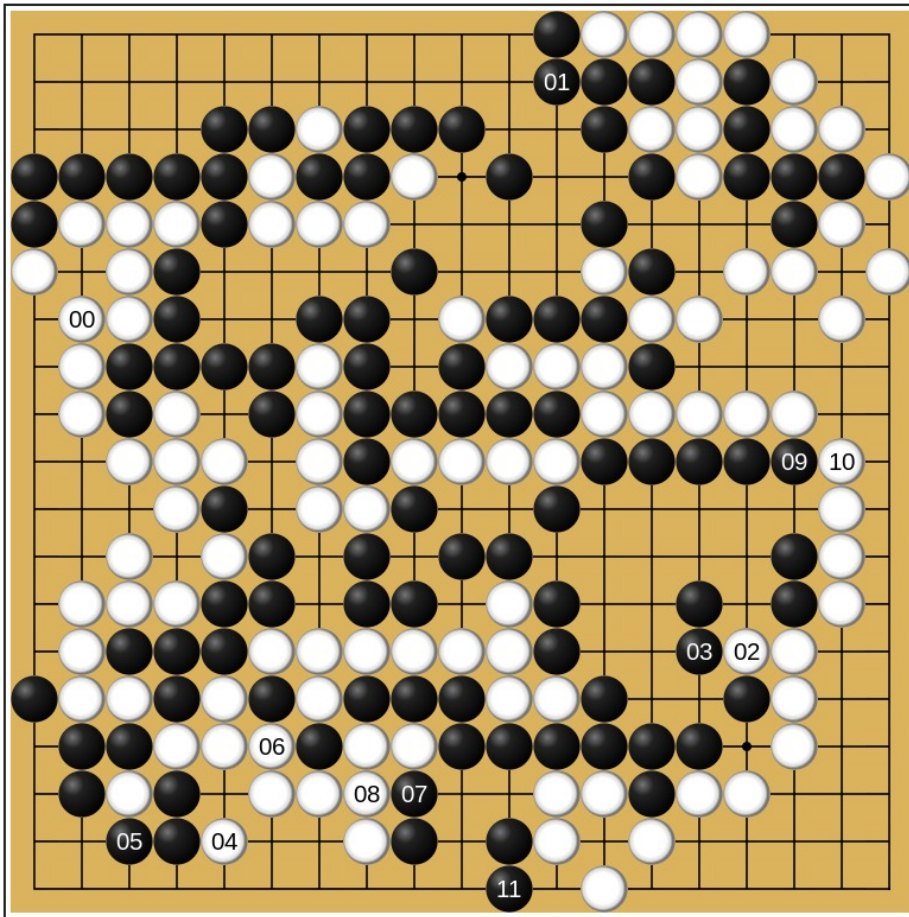
# Q – learning: convergence

- For Algorithm 3 (temporal difference learning),  $Q$  converges to  $Q^*$  if
  1. Every valid state-action pair is visited infinitely often
    - Q-learning is exploration-insensitive: any visitation strategy that satisfies this property will work!
  2.  $0 \leq \gamma < 1$
  3.  $\exists \beta$  s.t.  $|R(s, a)| < \beta \forall s \in \mathcal{S}, a \in \mathcal{A}$
  4. Initial  $Q$  values are finite
  5. Learning rate  $\alpha_t$  follows some “schedule” s.t.  
 $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$  e.g.,  $\alpha_t = 1/t+1$

# Deep Q-learning

- What if state-action spaces are continuous?
- Use a parametric function,  $Q(s, a; \Theta)$ , to approximate  $Q^*(s, a)$ 
  - Learn the parameters using SGD
  - Training data  $(\mathbf{s}_t, a_t, r_t, \mathbf{s}_{t+1})$  gathered online by the agent/learning algorithm
- If the approximator is a deep neural network => deep Q-learning

# AlphaGo (Black) vs. Lee Sedol (White) Game 2 final position (AlphaGo wins)



## Playing Go

19-by-19 board

Players alternate placing black and white stones

The goal is claim more territory than the opponent

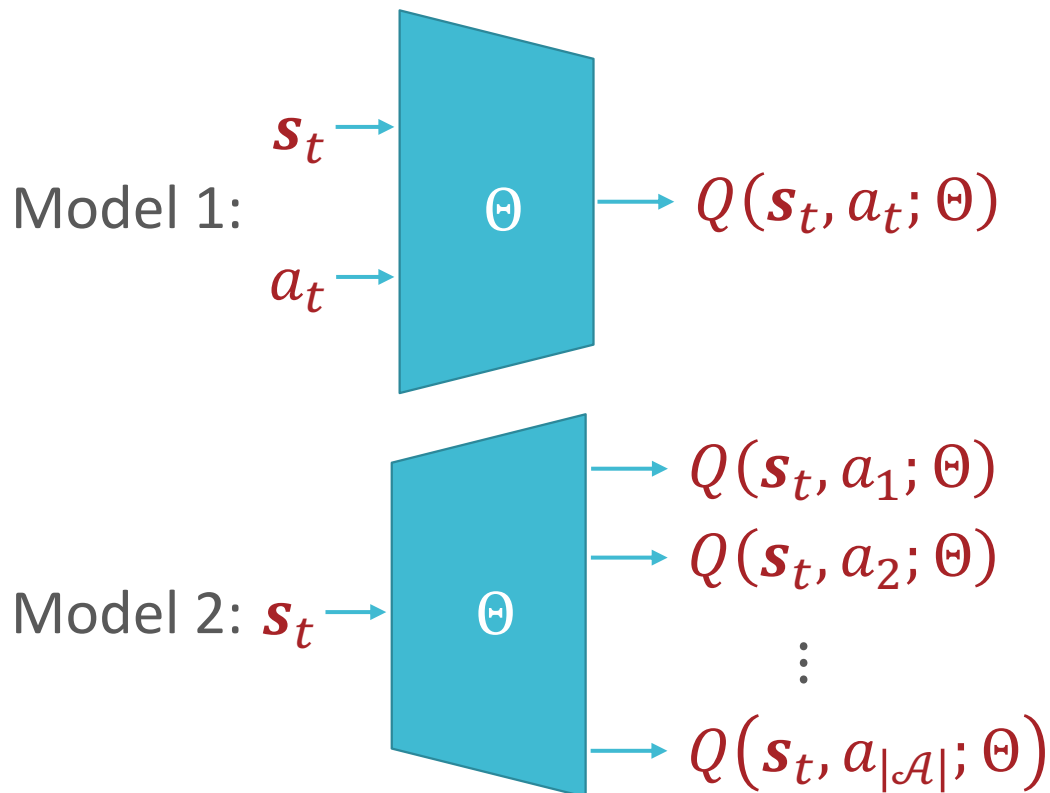
There are  $\sim 10^{170}$  legal Go board states!

Source: [https://en.wikipedia.org/wiki/AlphaGo\\_vs.\\_Lee\\_Sedol](https://en.wikipedia.org/wiki/AlphaGo_vs._Lee_Sedol)

Source: [https://en.wikipedia.org/wiki/Go\\_and\\_mathematics](https://en.wikipedia.org/wiki/Go_and_mathematics)

# Deep Q-learning: Model

- Represent states using some feature vector  $\mathbf{s}_t \in \mathbb{R}^M$   
e.g. for Go,  $\mathbf{s}_t = [1, 0, -1, \dots, 1]^T$
- Define a neural network architecture



# Deep Q-learning: Loss function

- “True” loss

$$\ell(\Theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \left( \overbrace{Q^*(s, a)}^{2. \text{ Don't know } Q^*} - Q(s, a; \Theta) \right)^2$$

1.  $\mathcal{S}$  too big to compute this sum

1. Use stochastic gradient descent: just consider one state-action pair in each iteration

2. Use temporal difference learning:

- Given current parameters  $\Theta^{(t)}$  the temporal difference target is

$$Q^*(s, a) \approx r + \gamma \max_{a'} Q(s', a'; \Theta^{(t)}) := y$$

- Set the parameters in the next iteration  $\Theta^{(t+1)}$  such that  $Q(s, a; \Theta^{(t+1)}) \approx y$

$$\ell(\Theta^{(t)}, \Theta^{(t+1)}) = \left( y - Q(s, a; \Theta^{(t+1)}) \right)^2$$



# Deep Q-learning: parametric online learning

- Inputs: discount factor  $\gamma$ , an initial state  $s_0$ ,  
learning rate  $\alpha$
- Initialize parameters  $\Theta^{(0)}$
- For  $t = 0, 1, 2, \dots$ 
  - Gather training sample  $(s_t, a_t, r_t, s_{t+1})$ , compute  $y$
  - Update  $\Theta^{(t)}$  by taking a step opposite the gradient
$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \alpha \nabla_{\Theta^{(t+1)}} \ell(\Theta^{(t)}, \Theta^{(t+1)})$$

where

$$\begin{aligned} & \nabla_{\Theta^{(t+1)}} \ell(\Theta^{(t)}, \Theta^{(t+1)}) \\ &= 2 \left( y - Q(s, a; \Theta^{(t+1)}) \right) \nabla_{\Theta^{(t+1)}} Q(s, a; \Theta^{(t+1)}) \end{aligned}$$

# Deep Q-learning: Experience replay

- Issue: SGD assumes i.i.d. training samples but in RL, samples are *highly* correlated
- Idea: keep a “replay memory”  $\mathcal{D} = \{e_1, e_2, \dots, e_N\}$  of the  $N$  most recent experiences  $e_t = (\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  (Lin, 1992)
  - Also keeps the agent from “forgetting” about recent experiences
- Alternate between:
  1. Sampling some  $e_i$  uniformly at random from  $\mathcal{D}$  and applying a Q-learning update (repeat  $T$  times)
  2. Adding a new experience to  $\mathcal{D}$
- Can also sample experiences from  $\mathcal{D}$  according to some distribution that prioritizes experiences with high error (Schaul et al., 2016)

# RL summary

- States, actions, rewards
- Policy
- Value function, Q function
- Finding optimal policy:
  - value iteration
  - policy iteration
- Unknown reward and transition function:
  - Q learning (including temporal difference)
- Continuous states and actions:
  - parametric models, deep Q learning
  - Experience replay