



# 10-301/10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

# Q-Learning

+

# Deep RL

Matt Gormley, Henry Chai, Hoda Heidari

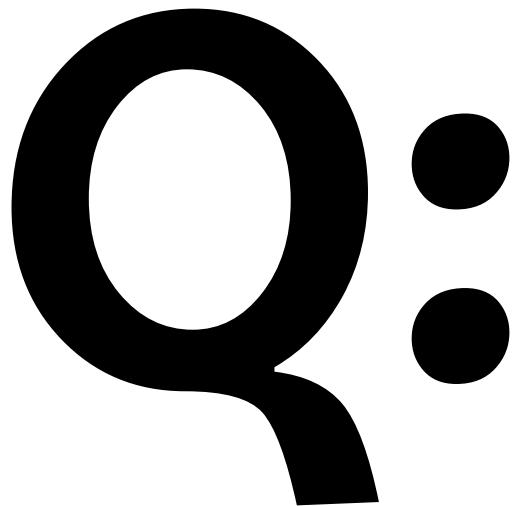
Lecture 22

Apr. 5, 2024

# Reminders

- **Homework 7: Deep Learning**
  - Out: Thu, Mar. 28
  - Due: Mon, Apr. 8 at 11:59pm
- **Schedule Notes**
  - Lecture 22: Fri, Apr. 5
  - HW8 Recitation: Mon, Apr. 8
- **Homework 8: Deep RL**
  - Out: Mon, Apr. 8
  - Due: Fri, Apr. 19 at 11:59pm

# **Q-LEARNING**



What can we do if we don't know the  
reward function / transition  
probabilities?

$$V(s) \leftarrow \max_a R(s,a) + \gamma \sum_{s'} p(s'|s,a) V(s')$$

Today's lecture is brought to you by the letter Q



Source: [https://en.wikipedia.org/wiki/Avenue\\_Q#/media/File:Image-AvenueQlogo.png](https://en.wikipedia.org/wiki/Avenue_Q#/media/File:Image-AvenueQlogo.png)

# Today's lecture is brought to you by the letter Q



Source: [https://vignette1.wikia.nocookie.net/jamesbond/images/9/9a/The\\_Four\\_Qs\\_-\\_Profile\\_\(2\).png/revision/latest?cb=20121102195112](https://vignette1.wikia.nocookie.net/jamesbond/images/9/9a/The_Four_Qs_-_Profile_(2).png/revision/latest?cb=20121102195112)

# Today's lecture is brought to you by the letter Q



Source: <https://www.npr.org/2017/06/03/531044118/there-may-not-be-flying-but-quidditch-still-creates-magic>

# Value Iteration

---

## Algorithm 1 Value Iteration

---

```
1: procedure VALUEITERATION( $R(s, a)$  reward function,  $p(\cdot|s, a)$ 
   transition probabilities)
2:   Initialize value function  $V(s) = 0$  or randomly
3:   while not converged do
4:     for  $s \in \mathcal{S}$  do
5:       for  $a \in \mathcal{A}$  do
6:         
$$Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a)V(s')$$

7:       
$$V(s) = \max_a Q(s, a)$$

8:     Let  $\pi(s) = \operatorname{argmax}_a Q(s, a)$ ,  $\forall s$ 
9:   return  $\pi$ 
```

---

Variant 1: with  $Q(s, a)$  table

# Q-Learning Motivation and $Q^*(s,a)$

$V^\pi(s)$

$Q^\pi(s,a)$  the expected discounted future reward of taking action  $a$  in state  $s$  and then following policy  $\pi$

- **Q-Learning Motivation**

Q: What if we don't know  $R(s,a)$  or  $p(s' | s, a)$ ?

A: Then value iteration and policy iteration don't work!

- **Definition:** Let  $Q^*(s,a)$  be the (true) expected discounted future reward of taking action  $a$  in state  $s$  and following optimal policy  $\pi^*$

$$V^*(s) = \max_a Q^*(s, a)$$

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

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

- **Key insight:** if we can learn  $Q^*$ , we can define  $\pi^*$  without knowing  $R(s,a)$  or  $p(s' | s, a)$ !

$$\pi^*(s) = \arg\max_{a \in A} Q^*(s, a)$$

non-deterministic version

# Q-Learning Motivation and $Q^*(s,a)$

- **Q-Learning Motivation**

Q: What if we don't know  $R(s,a)$  or  $\delta(s, a)$ ?

A: Then value iteration and policy iteration don't work!

- **Definition:** Let  $Q^*(s,a)$  be the (true) expected discounted future reward of taking action  $a$  in state  $s$

$$V^*(s) = \max_a Q^*(s, a)$$

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

$$Q^*(s, a) = R(s, a) + \gamma \left[ \max_{a'} Q^*(\underbrace{\delta(s, a), a'}_{s'}) \right]$$

- **Key insight:** if we can learn  $Q^*$ , we can define  $\pi^*$  without knowing  $R(s,a)$  or  $\delta(s, a)$ !

deterministic  
version

# Q-Learning Algorithm

deterministic  
version

produce training  
example  $(s, a, r, s')$

---

## Algorithm 1 Q-Learning (deterministic environment)

---

```
1: procedure QLEARNING( $\epsilon$ )
2:   Initialize  $Q$  function  $Q(s, a) = 0$  for all  $s, a$ 
3:   while true do
4:     select action  $a$  and execute
```

```
5:     receive reward  $r = R(s, a)$ 
6:     observe new state  $s' = \delta(s, a)$ 
7:     update table entry in  $Q$ 
```

we still don't know  $R$  or  $\delta$ ; these are given to agent by the environment

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

```
8:   Let  $\pi(s) = \operatorname{argmax}_a Q(s, a), \forall s$ 
9:   return  $\pi$ 
```

# Q-Learning Algorithm

deterministic  
version

produce training  
example  $(s, a, r, s')$

---

## Algorithm 1 Q-Learning (deterministic env., $\epsilon$ -greedy variant)

---

- 1: **procedure** QLEARNING( $\epsilon$ )  
   Initialize  $s$
- 2:    Initialize  $Q$  function  $Q(s, a) = 0$  for all  $s, a$
- 3:    **while** true **do**
- 4:     select action  $a$  and execute

exploit

with prob.  $(1 - \epsilon)$  : select  $a = \max_{a' \in \mathcal{A}} Q(s, a')$

with prob.  $\epsilon$  : select  $a \in \mathcal{A}$  randomly

explore

we still don't know  $R$  or  $\delta$ ; these are given to agent by the environment

- 5:     receive reward  $r = R(s, a)$
- 6:     observe new state  $s' = \delta(s, a)$
- 7:     update table entry in  $Q$

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

7.5       $s = s'$

- 8:     Let  $\pi(s) = \operatorname{argmax}_a Q(s, a), \forall s$
- 9:     **return**  $\pi$

# Q-Learning Algorithm

non-deterministic  
version

produce training  
example  $(s, a, r, s')$

$\alpha_n = \frac{1}{1 + \text{visits}(s, a, n)}$   
 $\text{visits}(s, a, n) = \# \text{ of visits to}$   
 $(s, a)$  up to and including step  $n$

---

## Algorithm 1 Q-Learning (non-deterministic env., $\epsilon$ -greedy variant)

---

1: **procedure** QLEARNING  
2:   Initialize  $Q$  function  $Q(s, a) = 0$  for all  $s, a$   
3:   **while** true **do**  
4:     select action  $a$  and execute

with prob.  $(1 - \epsilon)$  : select  $a = \max_{a' \in \mathcal{A}} Q(s, a')$

with prob.  $\epsilon$  : select  $a \in \mathcal{A}$  randomly

5:     receive reward  $r = R(s, a)$   
6:     observe new state  $s' = \cancel{\delta(s, a)} P(s' | s, a)$   
7:     update table entry in  $Q$

$$Q(s, a) \leftarrow (1 - \alpha_n)Q(s, a) + \alpha_n(r + \gamma \max_{a' \in \mathcal{A}} Q(s', a'))$$

Let  $\pi(s) = \operatorname{argmax}_a Q(s, a), \forall s$   
**return**  $\pi$

current value  
in the table

Q-learning update  
from deterministic  
version

# Q-Learning Algorithm

non-deterministic  
version

produce training  
example  $(s, a, r, s')$

$$\alpha_n = \frac{1}{1 + \text{visits}(s, a, n)}$$

$\text{visits}(s, a, n) = \#$  of visits to  $(s, a)$  up to and including step  $n$

---

## Algorithm 1 Q-Learning (non-deterministic env., $\epsilon$ -greedy variant)

---

```
1: procedure QLEARNING
2:   Initialize  $Q$  function  $Q(s, a) = 0$  for all  $s, a$ 
3:   while true do
4:     select action  $a$  and execute
       with prob.  $(1 - \epsilon)$  : select  $a = \max_{a' \in \mathcal{A}} Q(s, a')$ 
       with prob.  $\epsilon$  : select  $a \in \mathcal{A}$  randomly
5:     receive reward  $r = R(s, a)$ 
6:     observe new state  $s' = \delta(s, a)$ 
7:     update table entry in  $Q$ 
```

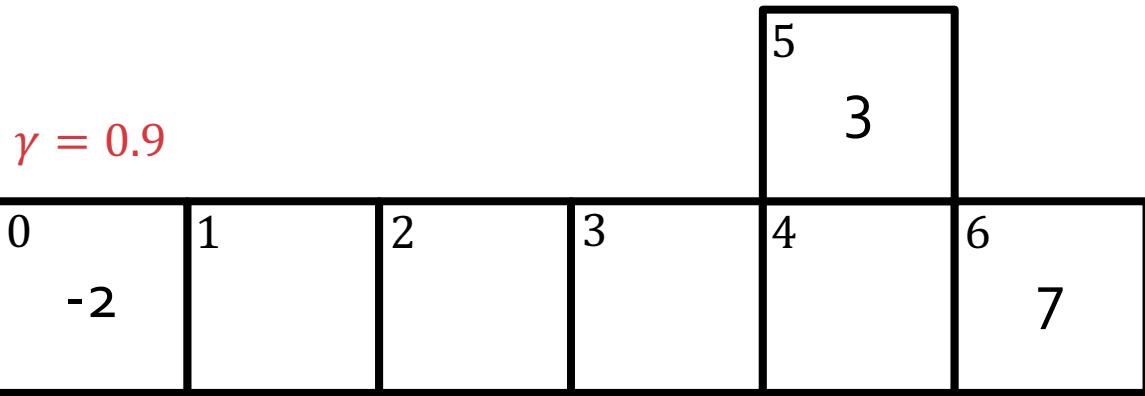
$$Q(s, a) \leftarrow Q(s, a) + \alpha_n(r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a))$$

Let  $\pi(s) = \operatorname{argmax}_a Q(s, a)$ ,  $\forall s$   
**return**  $\pi$

current value  
in the table

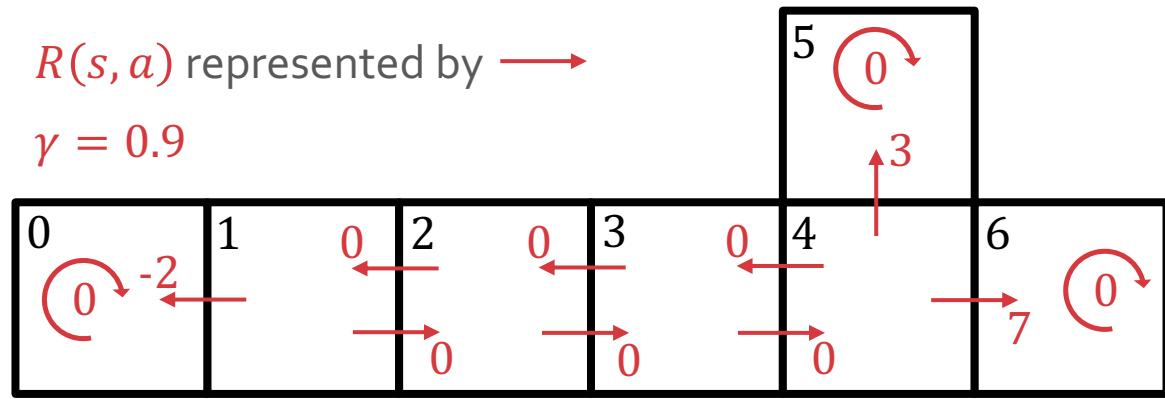
temporal difference  
target

# Learning $Q^*(s, a)$ : Example



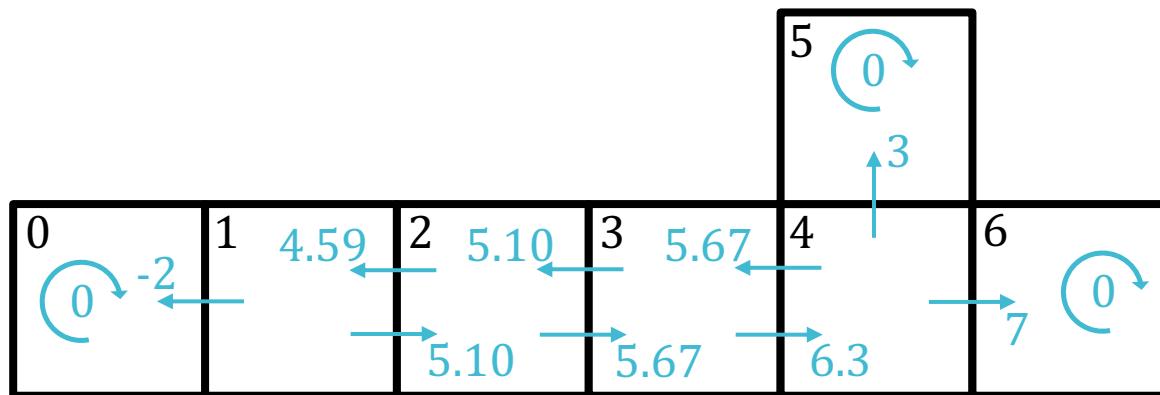
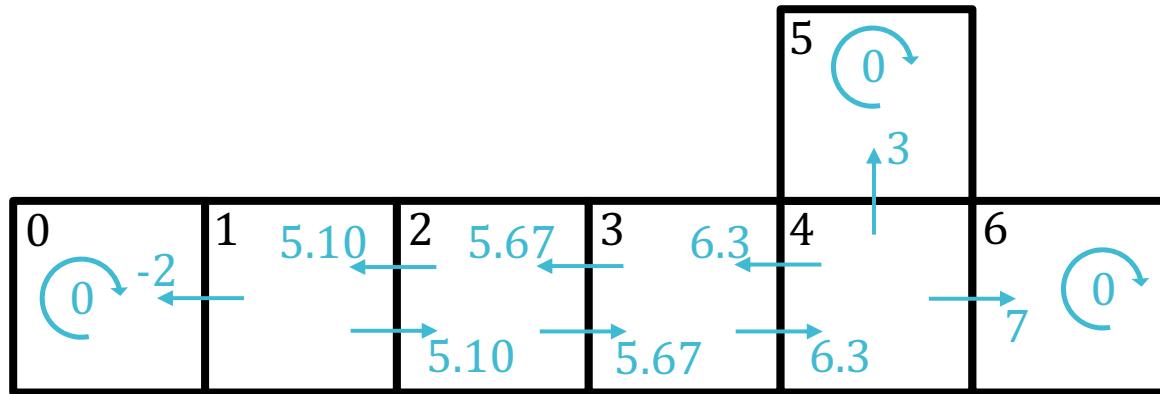
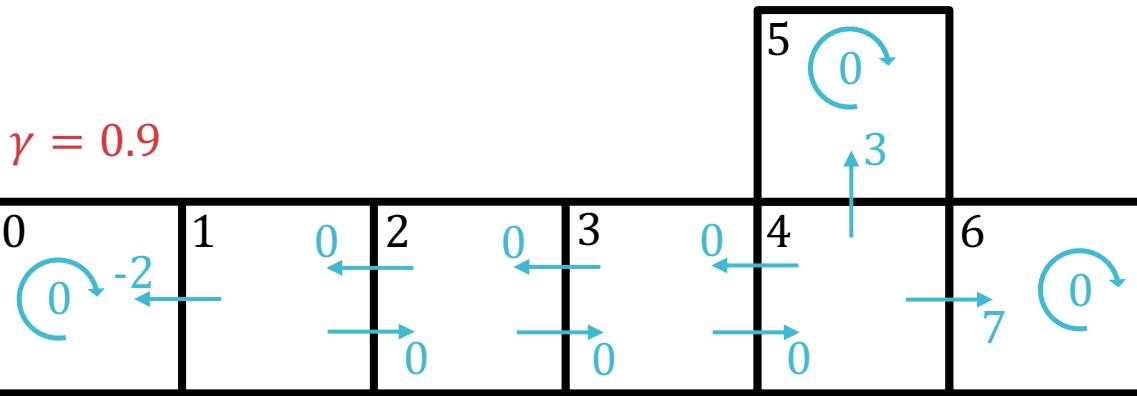
$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

# Learning $Q^*(s, a)$ : Example



Q1

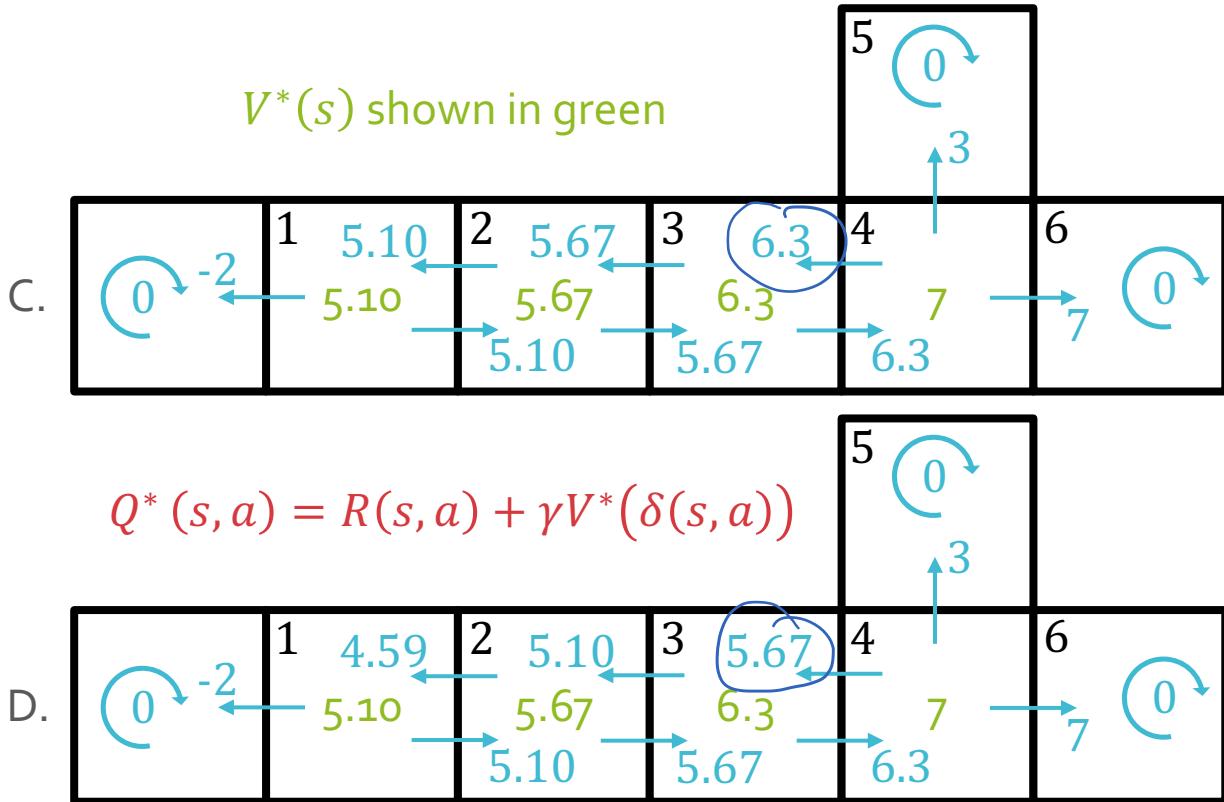
Poll: Which set of blue arrows (roughly) corresponds to  $Q^*(s, a)$ ?



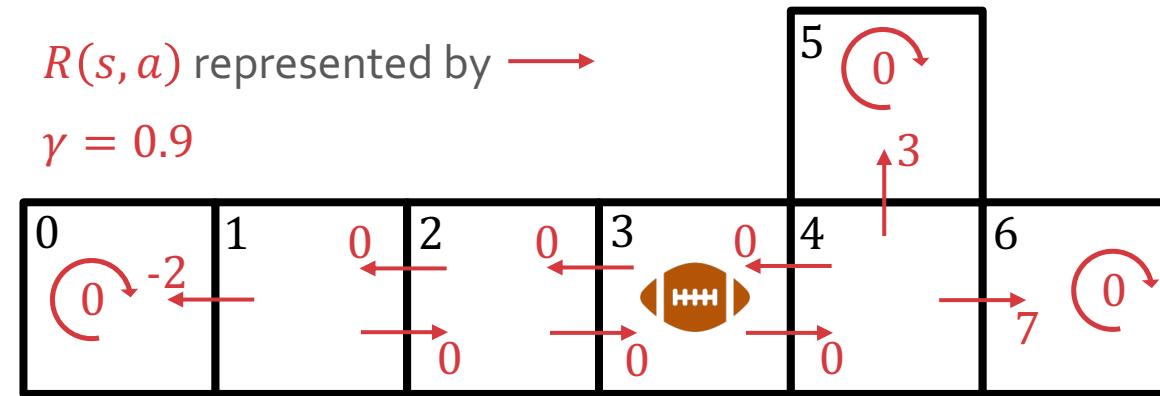
25%

60%

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

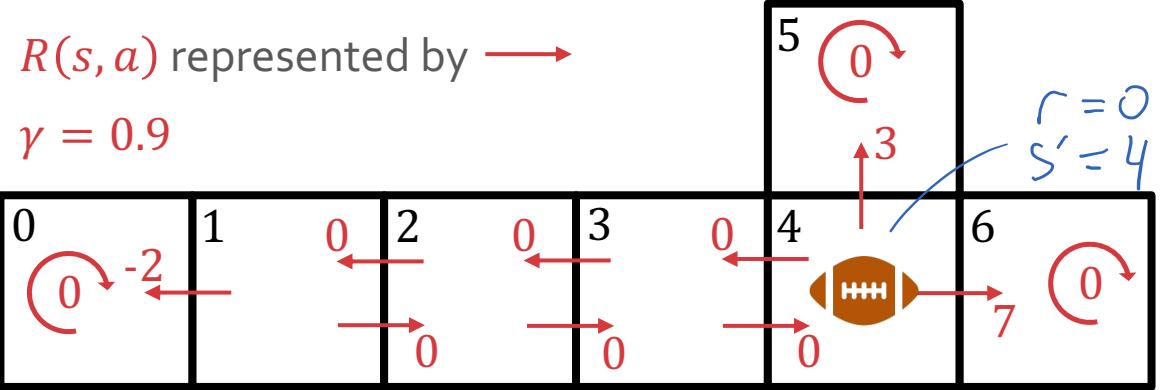


# Learning $Q^*(s, a)$ : Example



$Q(s, a)$	→	←	↑	↻
<b>0</b>	0	0	0	0
<b>1</b>	0	0	0	0
<b>2</b>	0	0	0	0
<b>3</b>	0	0	0	0
<b>4</b>	0	0	0	0
<b>5</b>	0	0	0	0
<b>6</b>	0	0	0	0

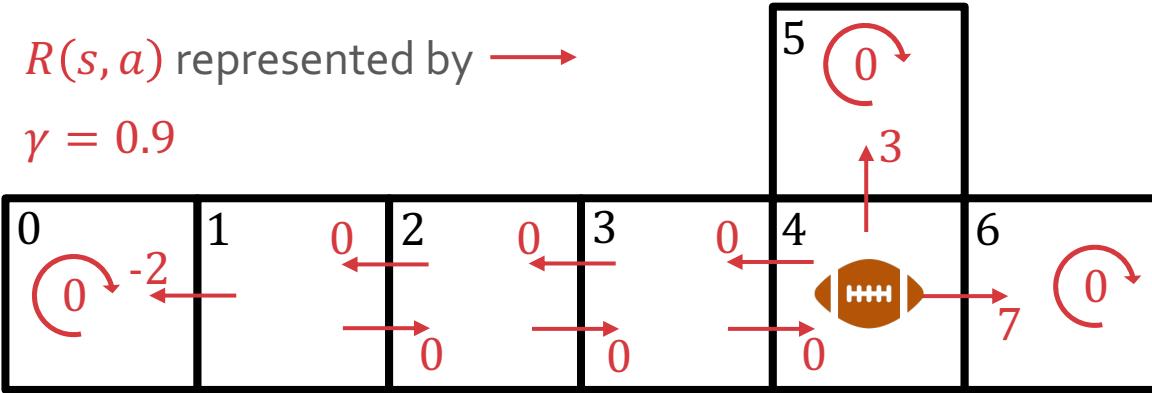
# Learning $Q^*(s, a)$ : Example



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

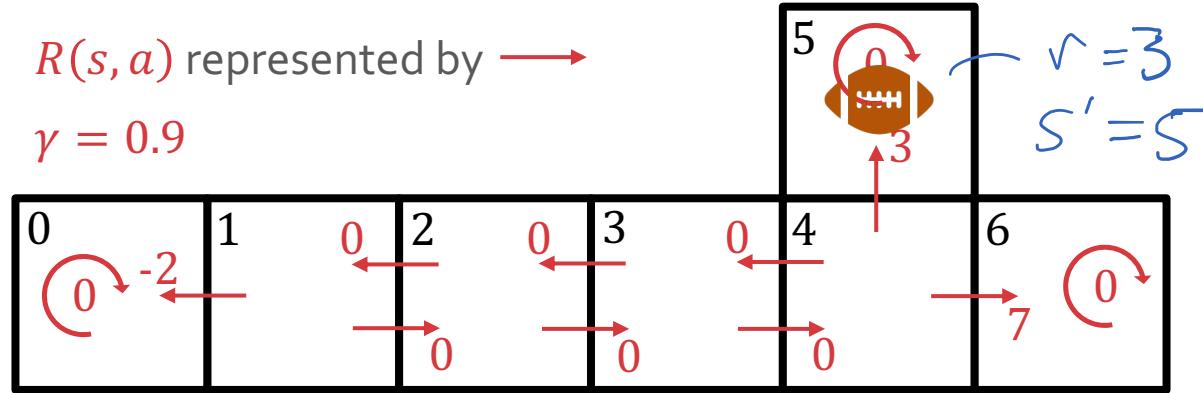
$Q(s, a)$	→	←	↑	↓
0	0	0	0	0
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0

# Learning $Q^*(s, a)$ : Example



$Q(s, a)$	→	←	↑	↻
<b>0</b>	0	0	0	0
<b>1</b>	0	0	0	0
<b>2</b>	0	0	0	0
<b>3</b>	0	0	0	0
<b>4</b>	0	0	0	0
<b>5</b>	0	0	0	0
<b>6</b>	0	0	0	0

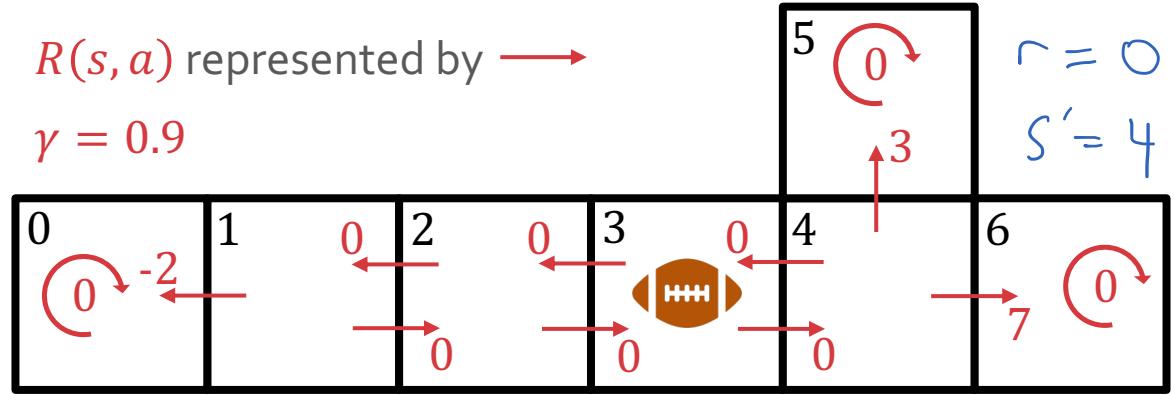
# Learning $Q^*(s, a)$ : Example



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

$Q(s, a)$	→	←	↑	↻
<b>0</b>	0	0	0	0
<b>1</b>	0	0	0	0
<b>2</b>	0	0	0	0
<b>3</b>	0	0	0	0
<b>4</b>	0	0	0	0
<b>5</b>	0	0	0	0
<b>6</b>	0	0	0	0

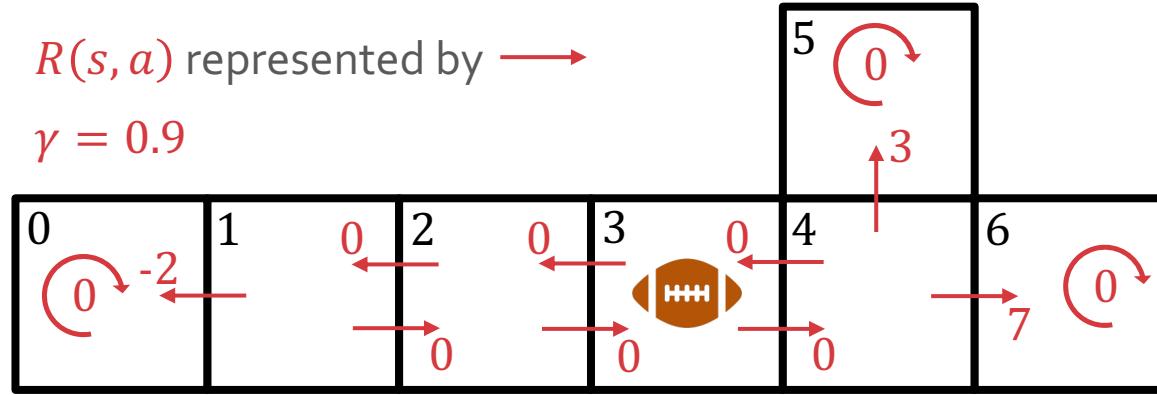
# Learning $Q^*(s, a)$ : Example



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

$Q(s, a)$	→	←	↑	↻
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

# Learning $Q^*(s, a)$ : Example



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

$Q(s, a)$	→	←	↑	↓
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

# Q-Learning Convergence

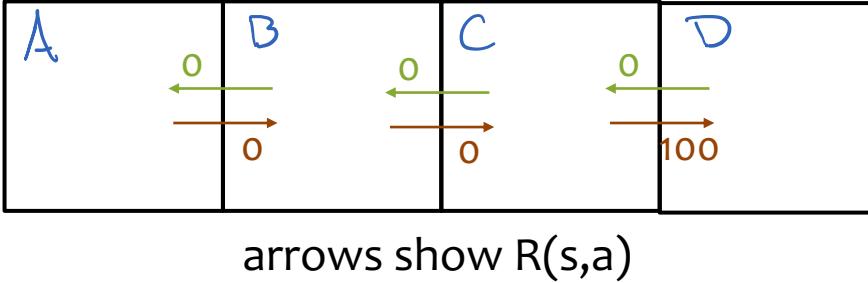
## Remarks

- Q converges to  $Q^*$  with probability 1.0, assuming...
  1. each  $\langle s, a \rangle$  is visited infinitely often
  2.  $0 \leq \gamma < 1$
  3. rewards are bounded  $|R(s,a)| < \beta$ , for all  $\langle s,a \rangle$
  4. initial Q values are finite
  5. Learning rate  $\alpha_t$  follows some “schedule” s.t.  $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 = 0$ , e.g.,  $\alpha_t = 1/(t+1)$
- Q-Learning is **exploration insensitive**  
⇒ visiting the states in any order will work assuming point 1 is satisfied
- May take **many iterations** to converge in practice



# Reordering Experiences

Easiest maze ever!



$$\begin{aligned}\gamma &= 0.9 \\ \mathcal{S} &= \{A, B, C, D\} \\ \mathcal{A} &= \{E, W\} \\ Q(s,a) &= 0 \text{ at the start}\end{aligned}$$

1. Suppose we visit states as below

i	s	a	r	s'
1	A	E	0	B
2	B	E	0	C
3	C	E	100	D

$$\begin{aligned}Q(A, E) &= 0 \\ Q(B, E) &= 0 \\ Q(C, E) &= 100\end{aligned}$$

2. Suppose we visit states **in reverse**

i	s	a	r	s'
1	C	E	100	D
2	B	E	0	C
3	A	E	0	B

$$\begin{aligned}Q(C, E) &= 100 \\ Q(B, E) &= 90 \\ Q(A, E) &= 81\end{aligned}$$

# Designing State Spaces

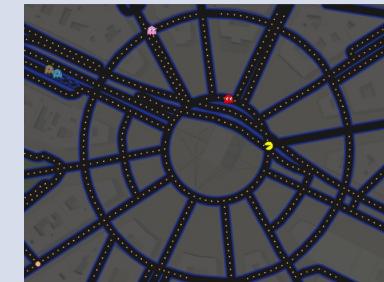
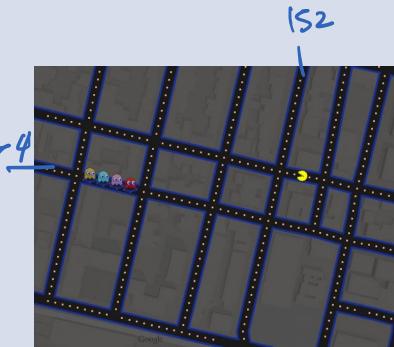
**Q:** Do we have to retrain our RL agent every time we change our state space?

**A:** Yes. But whether your state space changes from one setting to another is determined by your design of the state representation.

Two examples:

- State Space A:  $\langle x, y \rangle$  position on map  
e.g.  $s_t = \langle 74, 152 \rangle$
- State Space B: window of pixel colors centered at current Pac Man location  
e.g.  $s_t =$

0	1	0
0	0	0
1	1	1



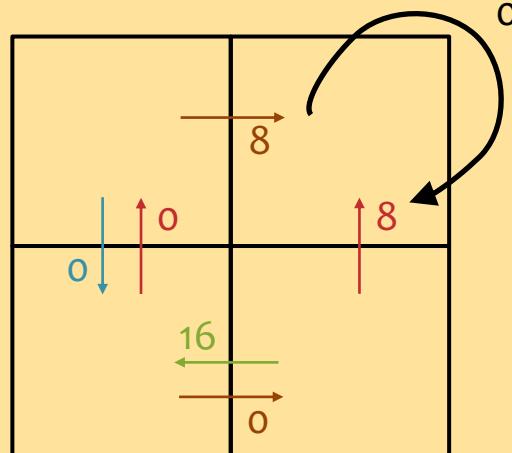
Q2 : skip

# Poll: Q-Learning

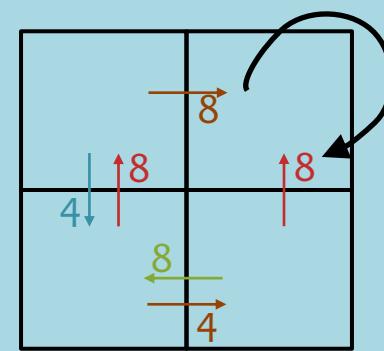
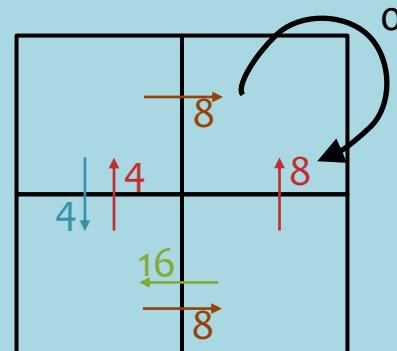
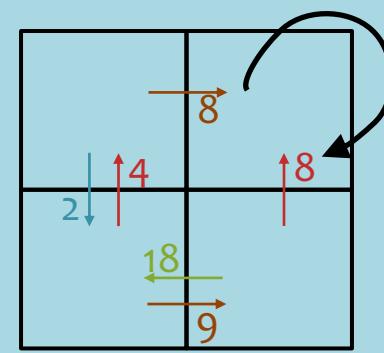
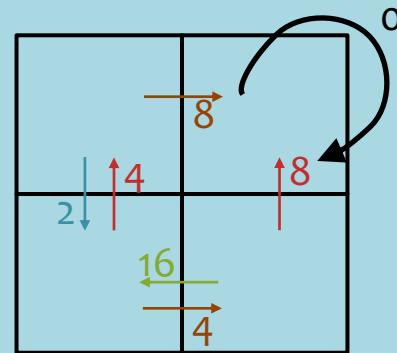
## Question:

For the  $R(s,a)$  values shown on the arrows below, which are the corresponding  $Q^*(s,a)$  values?

Assume discount factor = 0.5.



## Answer:

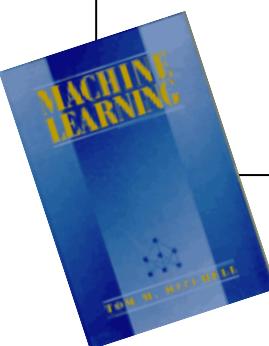
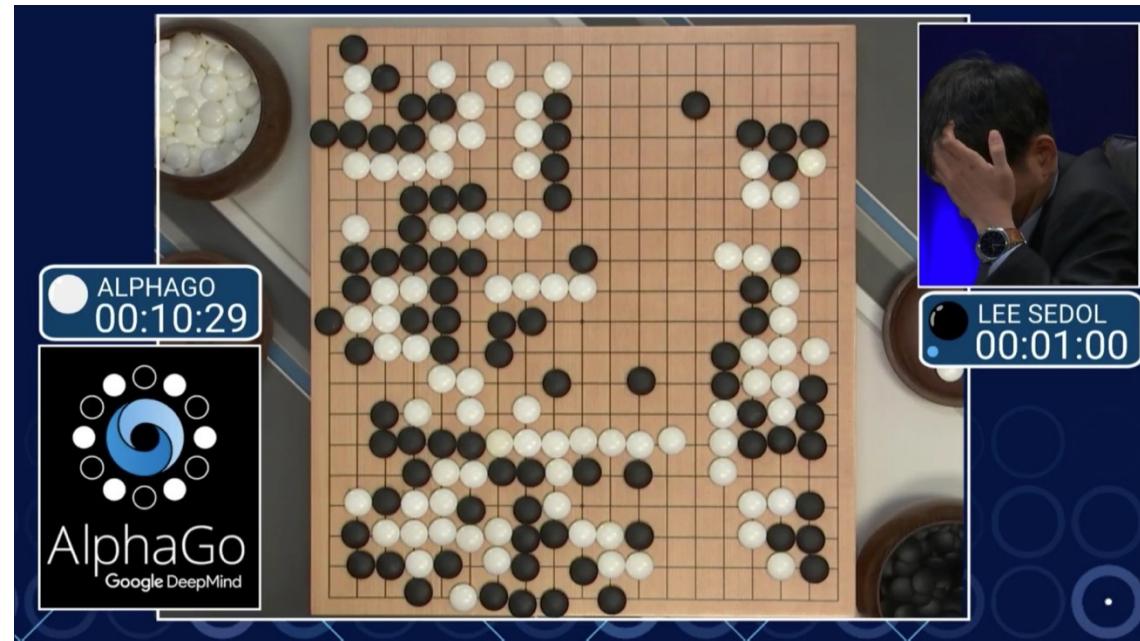


$F = \text{none } \nexists \text{ the above}$

# **DEEP RL FOR GAME OF GO**

# TD Gammon → Alpha Go

Learning to beat the masters at board games

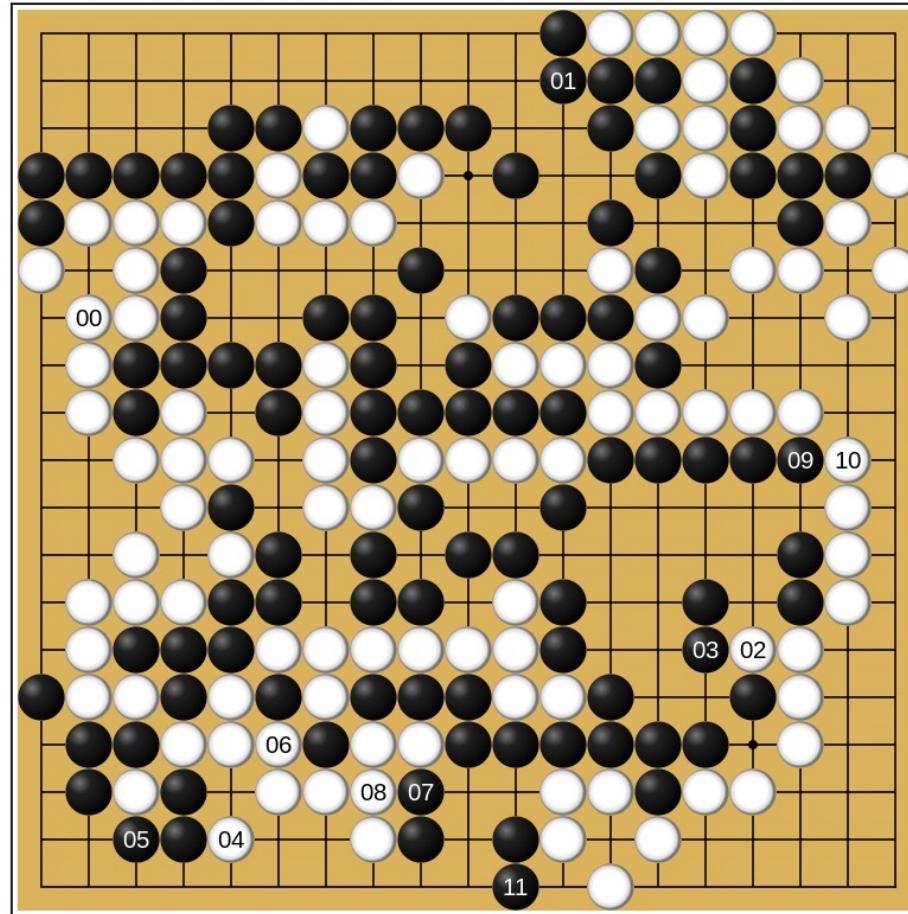
THEN	NOW
<p>“...the world’s top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”</p> <p>(Mitchell, 1997)</p> 	 <p>The image shows a Go board with black and white stones. On the left, there's a digital timer for AlphaGo showing 00:10:29. On the right, there's a digital timer for Lee Sedol showing 00:01:00. In the top right corner of the board area, there's a small inset video frame showing Lee Sedol, the South Korean Go player, with his hand near his face in a gesture of distress or exhaustion. The board itself is light-colored wood with a dark grid.</p>

# Alpha Go

## Game of Go (圍棋)

- 19x19 board
- Players alternately play black/white stones
- Goal is to fully encircle the largest region on the board
- Simple rules, but extremely complex game play

AlphaGo (Black) vs. Lee Sedol (White) - Game 2  
Final position (AlphaGo wins in 211 moves)



# Alpha Go

- State space is too large to represent explicitly since # of sequences of moves is  $O(b^d)$ 
  - Go:  $b=250$  and  $d=150$
  - Chess:  $b=35$  and  $d=80$
- Key idea:
  - Define a neural network to approximate the value function
  - Train by policy gradient

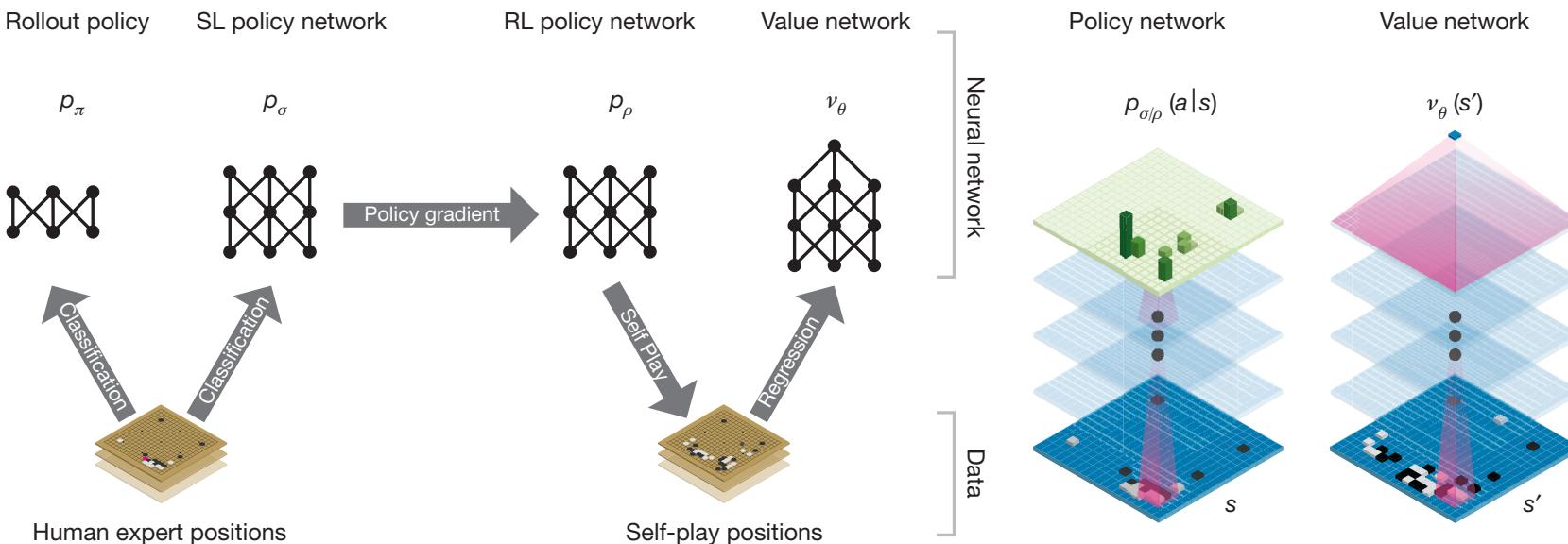
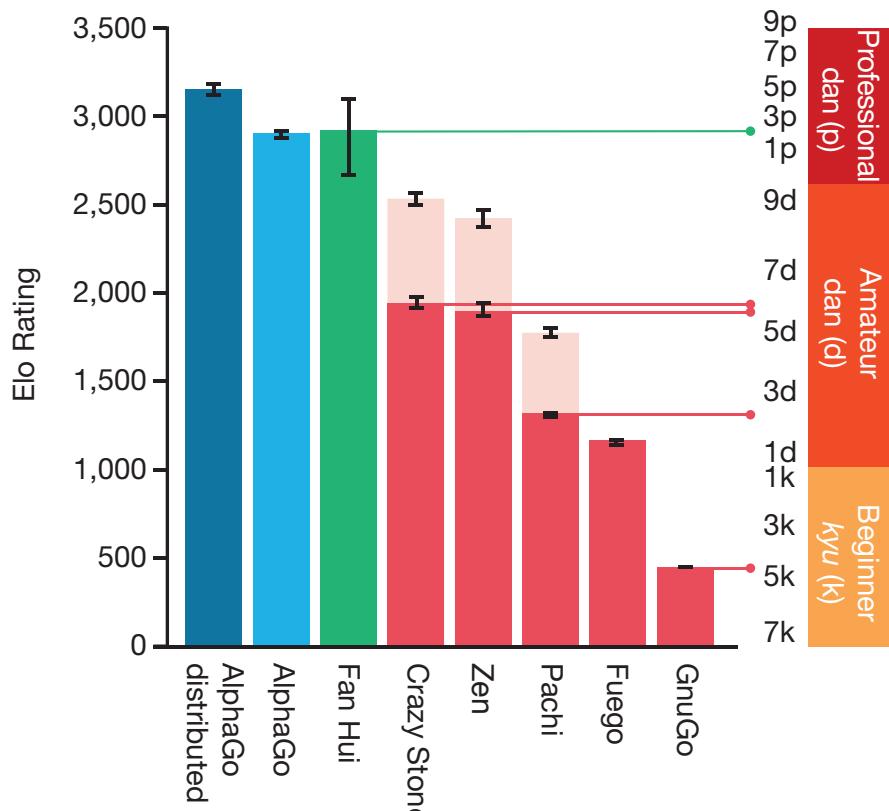


Figure from Silver et al. (2016)

# Alpha Go

- Results of a tournament
- From Silver et al. (2016): “a 230 point gap corresponds to a 79% probability of winning”



# **DEEP Q-LEARNING**

# Deep Q-Learning

**Question:** What if our state space  $S$  is too large to represent with a table?

**Examples:**

- $s_t$  = pixels of a video game
- $s_t$  = continuous values of sensors in a manufacturing robot
- $s_t$  = sensor output from a self-driving car

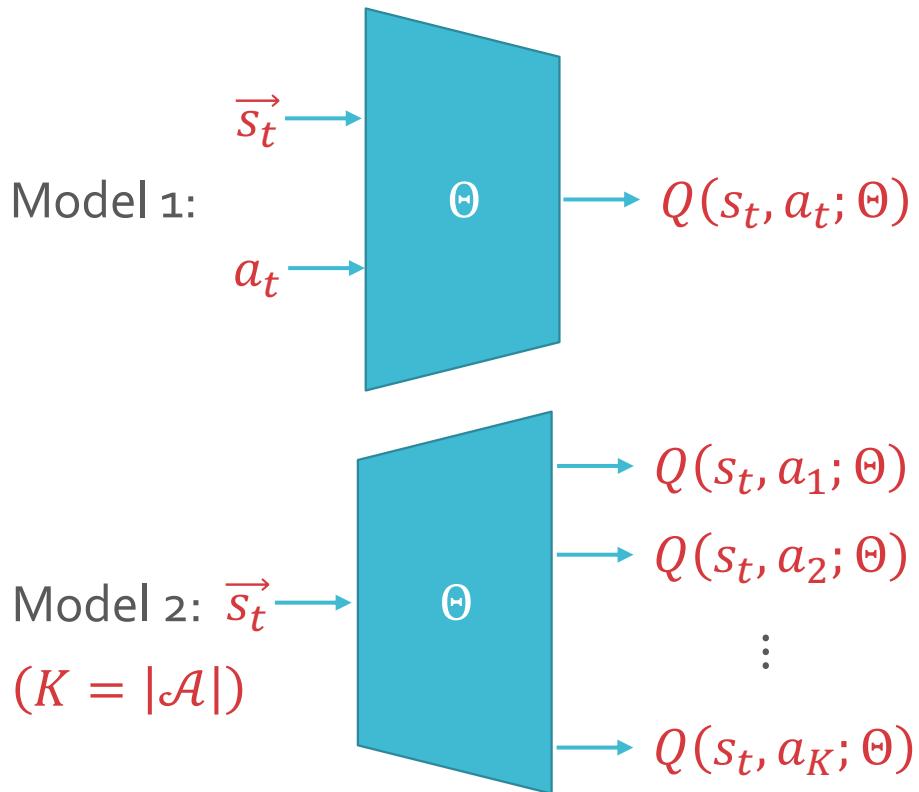
**Answer:** Use a parametric function to approximate the table entries

**Key Idea:**

1. Use a neural network  $Q(s,a; \theta)$  to approximate  $Q^*(s,a)$
2. Learn the parameters  $\theta$  via SGD with training examples  $\langle s_t, a_t, r_t, s_{t+1} \rangle$

# Deep Q-learning: Model

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



# Deep Q-learning: Model

- Represent states using some feature vector  $\vec{s}_t \in \mathbb{R}^M$   
e.g.,  $\vec{s}_t = [1, 0, 0, \dots, 1]^T$
- Define a neural network a bunch of linear regressors  
(technically still neural networks...), one for each action (let  $K = |\mathcal{A}|$ )

$$Q(\vec{s}, a_k; \Theta) = \vec{\theta}_k^T \vec{s} \text{ where } \Theta = \begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vdots \\ \vec{\theta}_K \end{bmatrix} \in \mathbb{R}^{K \times M}$$

- Goal:  $K \times M \ll |\mathcal{S}| \rightarrow$  computational tractability!
- Gradients are easy:  $\nabla_{\vec{\theta}_j} Q(\vec{s}, a_k; \Theta) = \begin{cases} \vec{0} & \text{if } j \neq k \\ \vec{s} & \text{if } j = k \end{cases}$

# Deep Q-learning: Model

- Represent states using some feature vector  $\vec{s}_t \in \mathbb{R}^M$   
e.g.,  $\vec{s}_t = [1, 0, 0, \dots, 1]^T$
- Define a neural network a bunch of linear regressors  
(technically still neural networks...), one for each action (let  $K = |\mathcal{A}|$ )

$$Q(\vec{s}, a_k; \Theta) = \vec{\theta}_k^T \vec{s} \text{ where } \Theta = \begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vdots \\ \vec{\theta}_K \end{bmatrix} \in \mathbb{R}^{K \times M}$$

- Goal:  $K \times M \ll |\mathcal{S}| \rightarrow$  computational tractability!

$$\nabla_{\Theta} Q(\vec{s}, a_k; \Theta) = \begin{bmatrix} \vec{0} \\ \vec{0} \\ \vdots \\ \vec{s} \\ \vdots \\ \vec{0} \end{bmatrix} \quad \leftarrow \text{Row } k$$

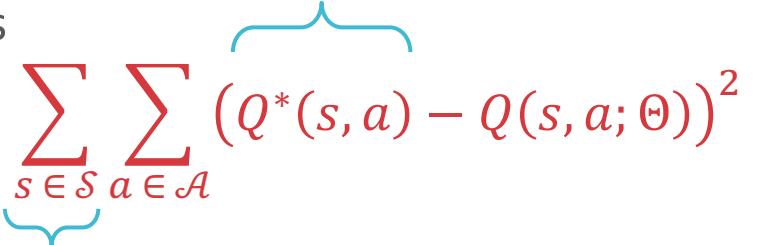
# Q-Learning (-ish) Update Rule

- Why don't we just do an update akin to what we do in regular Q-Learning?

Given  $(s, a, r, s')$ :

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

# Deep Q-learning: Loss Function

- “True” loss
- 
$$\ell(\Theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} (Q^*(s, a) - Q(s, a; \Theta))^2$$
- 2. Don’t know  $Q^*$
  - 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
$$Q^*(s, a) \approx r + \gamma \max_{a'} Q(s', a'; \Theta^{(t)}) \equiv y$$
      - Set the parameters in the next iteration  $\Theta^{(t+1)}$  such that  $Q(s, a; \Theta^{(t+1)}) \approx y$
- $$\ell(\Theta^{(t)}, \Theta^{(t+1)}) = (y - Q(s, a; \Theta^{(t+1)}))^2$$

# Deep Q-learning

- Algorithm: Online learning of  $Q^*$  (parametric form)
  - 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  $(s_t, a_t, r_t, s_{t+1})$
    - Update  $\Theta^{(t)}$  by taking a step opposite the gradient
    - $\Theta^{(const)} \leftarrow \Theta^{(t)}$   
$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \alpha \nabla_{\Theta^{(t)}} \ell(\Theta^{(const)}, \Theta^{(t)})$$

where

$$\begin{aligned}\nabla_{\Theta} \ell(\Theta^{(const)}, \Theta^{(t)}) &= 2 \left( y - Q(s, a; \Theta^{(t)}) \right) \nabla_{\Theta^{(t)}} Q(s, a; \Theta^{(t)}) \\ &= 2 \left( r + \gamma \max_{a'} Q(s', a'; \Theta^{(const)}) - Q(s, a; \Theta^{(t)}) \right) \nabla_{\Theta^{(t)}} Q(s, a; \Theta^{(t)})\end{aligned}$$

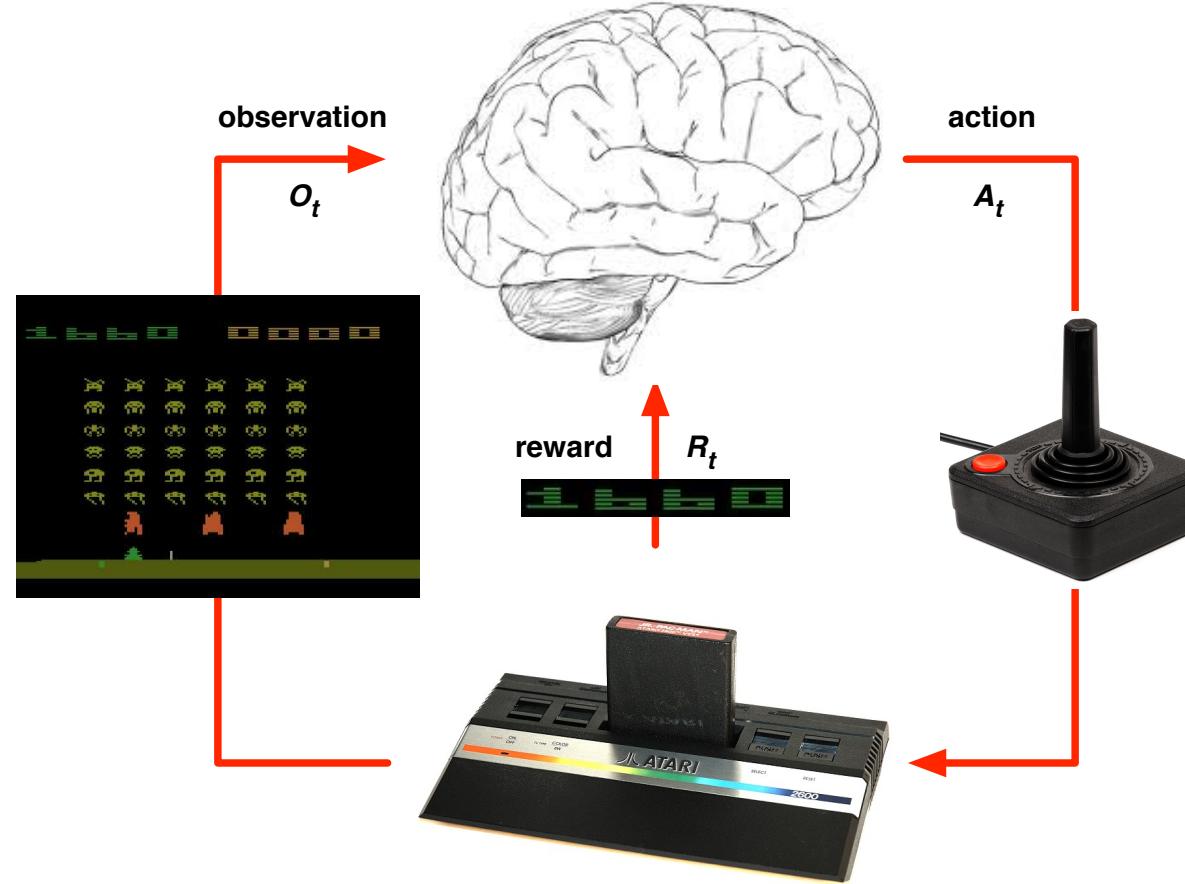
# Experience Replay

- **Problems** with online updates for Deep Q-learning:
  - not i.i.d. as SGD would assume
  - quickly forget rare experiences that might later be useful to learn from
- **Uniform Experience Replay** (Lin, 1992):
  - Keep a *replay memory*  $D = \{e_1, e_2, \dots, e_N\}$  of  $N$  most recent experiences  
 $e_t = \langle s_t, a_t, r_t, s_{t+1} \rangle$
  - Alternate two steps:
    1. Repeat  $T$  times: randomly sample  $e_i$  from  $D$  and apply a Q-Learning update to  $e_i$
    2. Agent selects an action using epsilon greedy policy to receive new experience that is added to  $D$
- **Prioritized Experience Replay** (Schaul et al, 2016)
  - similar to Uniform ER, but sample so as to prioritize experiences with high error

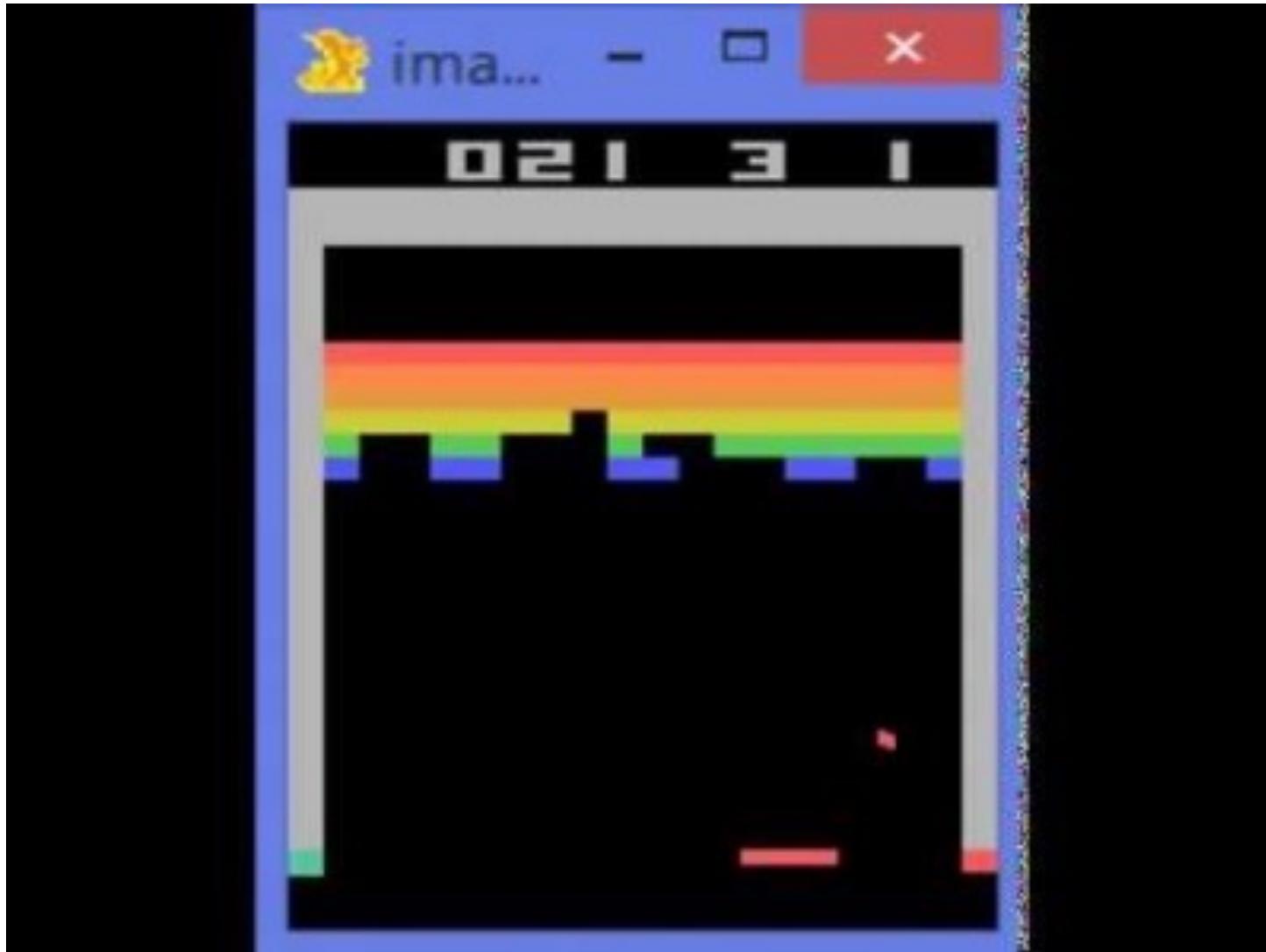
# **DEEP RL FOR ATARI GAMES**

# Playing Atari with Deep RL

- Setup: RL system observes the pixels on the screen
- It receives rewards as the game score
- Actions decide how to move the joystick / buttons



# Playing Atari games with Deep RL



Source: [https://www.youtube.com/watch?v=V1eYniJoRnk&t=2s&ab\\_channel=TwoMinutePapers](https://www.youtube.com/watch?v=V1eYniJoRnk&t=2s&ab_channel=TwoMinutePapers)

# Playing Atari games with Deep RL



Figure 1: Screen shots from five Atari 2600 Games: (*Left-to-right*) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

	B. Rider	Breakout	Enduro	Pong	Q*bert	Seaquest	S. Invaders
<b>Random</b>	354	1.2	0	-20.4	157	110	179
<b>Sarsa [3]</b>	996	5.2	129	-19	614	665	271
<b>Contingency [4]</b>	1743	6	159	-17	960	723	268
<b>DQN</b>	<b>4092</b>	<b>168</b>	<b>470</b>	<b>20</b>	<b>1952</b>	<b>1705</b>	<b>581</b>
<b>Human</b>	7456	31	368	-3	18900	28010	3690
<b>HNeat Best [8]</b>	3616	52	106	19	1800	920	<b>1720</b>
<b>HNeat Pixel [8]</b>	1332	4	91	-16	1325	800	1145
<b>DQN Best</b>	<b>5184</b>	<b>225</b>	<b>661</b>	<b>21</b>	<b>4500</b>	<b>1740</b>	1075

Table 1: The upper table compares average total reward for various learning methods by running an  $\epsilon$ -greedy policy with  $\epsilon = 0.05$  for a fixed number of steps. The lower table reports results of the single best performing episode for HNeat and DQN. HNeat produces deterministic policies that always get the same score while DQN used an  $\epsilon$ -greedy policy with  $\epsilon = 0.05$ .

# Learning Objectives

## **Reinforcement Learning: Q-Learning**

*You should be able to...*

1. Apply Q-Learning to a real-world environment
2. Implement Q-learning
3. Identify the conditions under which the Q-learning algorithm will converge to the true value function
4. Adapt Q-learning to Deep Q-learning by employing a neural network approximation to the Q function
5. Describe the connection between Deep Q-Learning and regression

# **BIG PICTURE**

# ML Big Picture

## Learning Paradigms:

*What data is available and when? What form of prediction?*

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

## Theoretical Foundations:

*What principles guide learning?*

- probabilistic
- information theoretic
- evolutionary search
- ML as optimization

## Problem Formulation:

*What is the structure of our output prediction?*

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

## Application Areas

*Key challenges?  
NLP, Speech, Computer Vision, Robotics, Medicine, Search*

## Facets of Building ML Systems:

*How to build systems that are robust, efficient, adaptive, effective?*

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

## Big Ideas in ML:

*Which are the ideas driving development of the field?*

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

# Learning Paradigms

Paradigm	Data
Supervised	$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \quad \mathbf{x} \sim p^*(\cdot) \text{ and } y = c^*(\cdot)$
↪ Regression	$y^{(i)} \in \mathbb{R}$
↪ Classification	$y^{(i)} \in \{1, \dots, K\}$
↪ Binary classification	$y^{(i)} \in \{+1, -1\}$
↪ Structured Prediction	$\mathbf{y}^{(i)}$ is a vector
Unsupervised	$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N \quad \mathbf{x} \sim p^*(\cdot)$
Semi-supervised	$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$
Online	$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \dots\}$
Active Learning	$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ and can query $y^{(i)} = c^*(\cdot)$ at a cost
Imitation Learning	$\mathcal{D} = \{(s^{(1)}, a^{(1)}), (s^{(2)}, a^{(2)}), \dots\}$
Reinforcement Learning	$\mathcal{D} = \{(s^{(1)}, a^{(1)}, r^{(1)}), (s^{(2)}, a^{(2)}, r^{(2)}), \dots\}$