

<https://introml.mit.edu/>

# 6.390 Intro to Machine Learning

## Lecture 11: Reinforcement Learning

Shen Shen

April 26, 2024

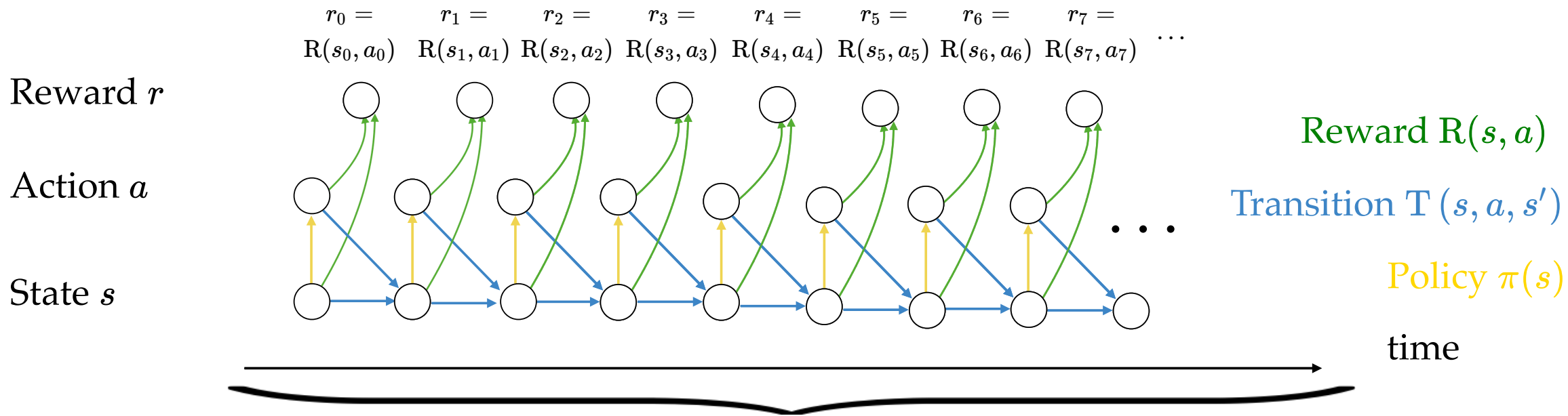
# Outline

- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?

# MDP Definition and Goal

- $\mathcal{S}$  : state space, contains all possible states  $s$ .
- $\mathcal{A}$  : action space, contains all possible actions  $a$ .
- $T(s, a, s')$  : the probability of transition from state  $s$  to  $s'$  when action  $a$  is taken.
- $R(s, a)$  : a function that takes in the (state, action) and returns a reward.
- $\gamma \in [0, 1]$ : discount factor, a scalar.
  
- $\pi(s)$  : policy, takes in a state and returns an action.

Ultimate goal of an MDP: Find the "best" policy  $\pi$ .



a trajectory (aka an experience or rollout)  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$

how "good" is a trajectory?

$$R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \gamma^3 R(s_3, a_3) + \gamma^4 R(s_4, a_4) + \gamma^5 R(s_5, a_5) + \gamma^6 R(s_6, a_6) + \gamma^7 R(s_7, a_7) \dots$$



## Running example: Mario in a grid-world

1	2	3
4	5	6
7	8	9

- 9 possible **states**
- 4 possible **actions**: {Up  $\uparrow$ , Down  $\downarrow$ , Left  $\leftarrow$ , Right  $\rightarrow$ }
- almost all **transitions** are deterministic:
  - Normally, actions take Mario to the “intended” state.
    - E.g., in state (7), action “ $\uparrow$ ” gets to state (4)
  - If an action would've taken us out of this world, stay put
    - E.g., in state (9), action “ $\rightarrow$ ” gets back to state (9)
  - except, in state (6), action “ $\uparrow$ ” leads to two possibilities:
    - 20% chance ends in (2)
    - 80% chance ends in (3)





# example cont'd

1	2	3
4	5	6
7	8	9

Transitions from state 2: Up (80%), Down (20%), Left, Right

		1
		1 1
		1
		-10
		-10 -10
		-10

- (state, action) pair can get Mario rewards:

- In state (3), any action gets reward +1 
- In state (6), any action gets reward -10 
- Any other (state, action) pairs get reward 0

actions: {Up ↑, Down ↓,  
Left ←, Right →}

- goal is to find a gameplay **policy** strategy for Mario, to get maximum expected sum of discounted rewards, with a **discount facotor**  $\gamma = 0.9$

Recall:

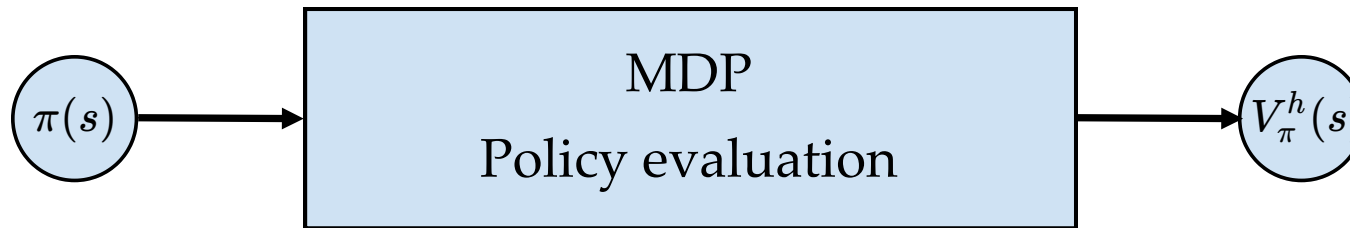
1	2	3
4	5	6
7	8	9

$\pi(s) = \text{“}\uparrow\text{”}, \forall s$

$R(3, \uparrow) = 1$

$R(6, \uparrow) = -10$

$\gamma = 0.9$


 $V_{\pi}^0(s)$ 

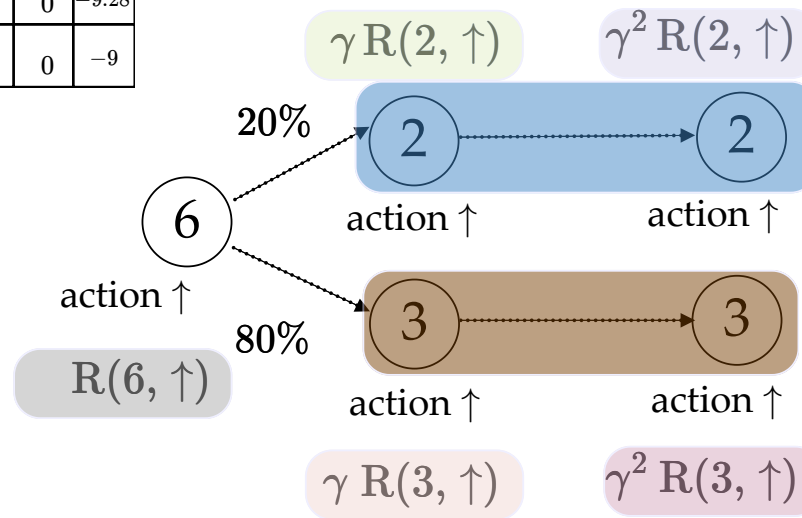
0	0	0
0	0	0
0	0	0

 $V_{\pi}^1(s)$ 

0	0	1
0	0	-10
0	0	0

 $V_{\pi}^2(s)$ 

0	0	1.9
0	0	-9.28
0	0	-9



Now, let's think about  $V_{\pi}^3(6)$

$$\begin{aligned}
 V_{\pi}^3(6) &= R(6, \uparrow) + 20\% [\gamma R(2, \uparrow) + \gamma^2 R(2, \uparrow)] + 80\% [\gamma R(3, \uparrow) + \gamma^2 R(3, \uparrow)] \\
 &= R(6, \uparrow) + 20\% \gamma [R(2, \uparrow) + \gamma R(2, \uparrow)] + 80\% \gamma [R(3, \uparrow) + \gamma R(3, \uparrow)] \\
 &= R(6, \uparrow) + 20\% \gamma V_{\pi}^2(2) + 80\% \gamma V_{\pi}^2(3)
 \end{aligned}$$

## finite-horizon policy evaluation

For a given policy  $\pi(s)$ , the finite-horizon horizon- $h$  (state) value functions are:

$$V_{\pi}^h(s) := \mathbb{E} \left[ \sum_{t=0}^{h-1} \gamma^t \mathbf{R}(s_t, \pi(s_t)) \mid s_0 = s, \pi \right], \forall s$$

### Bellman recursion

$$V_{\pi}^h(s) = \mathbf{R}(s, \pi(s)) + \gamma \sum_{s'} \mathbf{T}(s, \pi(s), s') V_{\pi}^{h-1}(s'), \forall s$$

## infinite-horizon policy evaluation

For any given policy  $\pi(s)$ , the infinite-horizon (state) value functions are

$$V_{\pi}(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{R}(s_t, \pi(s_t)) \mid s_0 = s, \pi \right], \forall s$$

$\gamma$  is now necessarily  $< 1$  for convergence too in general

### Bellman equation

$$V_{\pi}(s) = \mathbf{R}(s, \pi(s)) + \gamma \sum_{s'} \mathbf{T}(s, \pi(s), s') V_{\pi}(s'), \forall s$$

- $|\mathcal{S}|$  many linear equations





example: recursively finding  $Q^h(s, a)$

$Q^h(s, a)$  is the expected sum of discounted rewards for

- starting in state  $s$ ,
- take action  $a$ , for one step
- act **optimally** there afterwards for the remaining  $(h - 1)$  steps

$Q^0(s, a)$

0	0	0
0	0	0
0	0	0
0	0	0
0	0	0

$Q^1(s, a)$

0	0	1
0	0	1
0	0	-10
0	0	-10
0	0	-10
0	0	0
0	0	0
0	0	0

Recall:  $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Transitions: 2 to 3 (80%), 2 to 5 (20%), 5 to 6 (100%)

$R(s, a)$

		1
		1
		-10
		-10
		-10



$Q^h(s, a)$  is the expected sum of discounted rewards for

- starting in state  $s$ ,
- take action  $a$ , for one step
- act **optimally** there afterwards for the remaining  $(h - 1)$  steps

Recall:  $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states. State 2 is highlighted in yellow and state 3 is highlighted in purple. A dashed arrow points from state 6 to state 2, labeled '20%'. A solid arrow points from state 6 to state 3, labeled '80%'.

$Q^1(s, a) = R(s, a)$

0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

Diagram showing a 3x3 grid of states. The top-right cell (state 3) is highlighted in purple and contains the value 1. The middle-right cell (state 6) is highlighted in green and contains the value -10. The bottom-right cell (state 9) is highlighted in purple and contains the value -10.

$Q^2(s, a)$

					1.9
				1	1.9
					-8
					-9.28

Diagram showing a 3x3 grid of states. The top-right cell (state 3) is highlighted in purple and contains the value 1.9. The middle-right cell (state 6) is highlighted in green and contains the value -9.28.

- receive  $R(6, \uparrow)$
- act **optimally** for one more timestep, at the next state  $s'$ 
  - 20% chance,  $s' = 2$ , act optimally, receive  $\max_{a'} Q^1(2, a')$
  - 80% chance,  $s' = 3$ , act optimally, receive  $\max_{a'} Q^1(3, a')$

Let's consider  $Q^2(6, \uparrow) = R(6, \uparrow) + \gamma[.2 \max_{a'} Q^1(2, a') + .8 \max_{a'} Q^1(3, a')]$

$$= -10 + .9[.2 * 0 + .8 * 1] = -9.28$$



$Q^h(s, a)$  is the expected sum of discounted rewards for

- starting in state  $s$ ,
- take action  $a$ , for one step
- act **optimally** there afterwards for the remaining  $(h - 1)$  steps

Recall:  $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states. State 2 is highlighted in yellow. A transition from state 6 to state 2 is shown with a probability of 20% (indicated by a dashed arrow) and a transition from state 6 to state 3 is shown with a probability of 80% (indicated by a solid arrow).

$Q^1(s, a)$   
=  $R(s, a)$

0	0	0	0	1	1
0	0	0	0	1	1
0	0	0	0	-10	-10
0	0	0	0	-10	-10
0	0	0	0	0	0
0	0	0	0	0	0

Diagram showing a 3x3 grid of states. The top-right cell (state 3) is highlighted in purple. The middle-left cell (state 5) is highlighted in yellow.

$Q^2(s, a)$

				1.9	
				1	1.9
				-8	
				-9.28	

Diagram showing a 3x3 grid of states. The bottom-right cell (state 9) is highlighted in green.

$$Q^2(6, \uparrow) = R(6, \uparrow) + \gamma[.2 \max_{a'} Q^1(2, a') + .8 \max_{a'} Q^1(3, a')]$$

in general  $Q^h(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a'), \forall s, a$



$Q^h(s, a)$  is the expected sum of discounted rewards for

- starting in state  $s$ ,
- take action  $a$ , for one step
- act **optimally** there afterwards for the remaining  $(h - 1)$  steps

Recall:  $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states. State 2 is at the top-middle, state 3 is at the top-right, and state 5 is at the middle-middle. A dashed arrow points from state 2 to state 5, labeled '20%'. A solid arrow points from state 3 to state 5, labeled '80%'.

$Q^1(s, a)$

0	0	0	1	1
0	0	0	1	1
0	0	0	-10	-10
0	0	0	-10	-10
0	0	0	0	0
0	0	0	0	0

$Q^2(s, a)$

		1.9
		1 1.9
		-8
		-9.28

what's the optimal action in state 3, with horizon 2, given by  $\pi_2^*(3) = ?$

either up or right

in general

$$\pi_h^*(s) = \arg \max_a Q^h(s, a), \forall s, h$$

Given the finite horizon recursion

$$Q^h(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$$

We should easily be convinced of the infinite horizon equation

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

Infinite-horizon Value Iteration

1. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :

2.  $Q_{\text{old}}(s, a) = 0$

3. **while** True:

4. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :

5.  $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

6. **if**  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :

7. **return**  $Q_{\text{new}}$

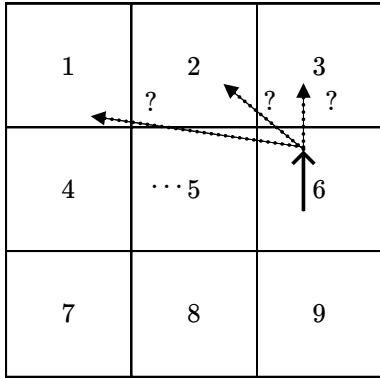
8.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

# Outline

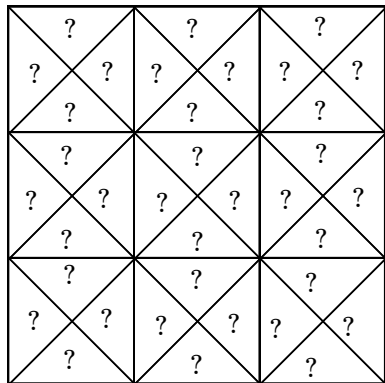
- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?



## Running example: Mario in a grid-world (the Reinforcement-Learning Setup)



- 9 possible **states**
- 4 possible **actions**: {Up  $\uparrow$ , Down  $\downarrow$ , Left  $\leftarrow$ , Right  $\rightarrow$ }
- **all transitions** probabilities are **unknown**.



- (state, action) pair gets Mario **unknown rewards**.
- goal is to find a gameplay **policy** strategy for Mario, to get maximum expected sum of discounted rewards, with a **discount facotor**  $\gamma = 0.9$

# RL Definition and Goal

- $\mathcal{S}$  : state space, contains all possible states  $s$ .
- $\mathcal{A}$  : action space, contains all possible actions  $a$ .
- $T(s, a, s')$  : the probability of transition from state  $s$  to  $s'$  when action  $a$  is taken.
- $R(s, a)$  : a function that takes in the (state, action) and returns a reward
- $\gamma \in [0, 1]$ : discount factor, a scalar.
- $\pi(s)$  : policy, takes in a state and returns an action.

Ultimate goal of an RL: Find the "best" policy  $\pi$ .



# Outline

- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?

## (MDP)-Model-Based Methods (for solving RL)

Keep playing the game to approximate the unknown rewards and transitions.

Rewards are particularly easy:

e.g. by observing what reward  $r$  received from being in state 6 and take  $\uparrow$  action, we know  $R(6, \uparrow)$

Transitions are a bit more involved but still simple:

e.g. play the game 1000 times, count the # of times (we started in state 6, take  $\uparrow$  action, end in state 2), then, roughly,  $T(6, \uparrow, 2) = (\text{that count}/1000)$

Now, with  $R, T$  estimated, we're back in MDP setting.

In Reinforcement Learning:

- *Model* typically means MDP tuple  $\langle \mathcal{S}, \mathcal{A}, T, R, \gamma \rangle$
- The learning objective is not referred to as *hypothesis* explicitly, we simply just call it the *policy*.

# Outline

- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?

How do we learn a good policy without learning transition or rewards explicitly?

We kinda already know a way: Q functions!

## 2.4) (Recall from MDP lab)

We switch to an infinite-horizon scenario. For our stochastic machine, here is the infinite-horizon  $Q$  function (computed via value iteration) for  $\gamma$  near 1.

```
      wash      paint      eject
dirty  [[ 2.32274541 -0.70048204  0.      ]
clean  [ 2.32274541  5.71581775  0.      ]
painted [ 2.32274541  6.9        10.     ]
ejected [ 0.          0.          0.     ]]
```

What is the optimal thing to do with a **clean** object?

What will you do if it becomes **dirty**?

Does this optimal policy make intuitive sense?

So once we have "good" Q values, we can find optimal policy easily.

But didn't we calculate this Q-table via value iteration using transition and rewards explicitly?

Indeed, recall that, in MDP:

- Finite horizon recursion

$$Q^h(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$$

- Infinite horizon equation

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

- Infinite-horizon Value Iteration

1. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :

2.  $Q_{\text{old}}(s, a) = 0$

3. **while** True:

4. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :

5.  $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

6. **if**  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :

7. **return**  $Q_{\text{new}}$

8.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

- value iteration relied on having full access to R and T

$$Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$$

- hmm... perhaps, we could simulate  $(s, a)$ , observe  $r$  and  $s'$ , and just use

$$r + \gamma \max_{a'} Q_{\text{old}}(s', a')$$

as the proxy for the r.h.s. assignment?

- BUT, this is basically saying the realized  $s'$  is the only possible next state; pretty rough! We'd override any previous "learned" Q values.

e.g.  $Q_{\text{new}}(6, \uparrow) \leftarrow -10 + \gamma \max_{a'} Q_{\text{old}}(2, a')$

$Q_{\text{new}}(6, \uparrow) \leftarrow -10 + \gamma \max_{a'} Q_{\text{old}}(3, a')$

- **better** way is to smoothly keep track of what's our old belief with new evidence:

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha) Q_{\text{old}}(s, a) + \alpha \left( r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

old belief

learning rate

target

## VALUE-ITERATION ( $\mathcal{S}, \mathcal{A}, T, R, \gamma, \epsilon$ )

1. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :
2.  $Q_{\text{old}}(s, a) = 0$
3. **while** True:
4. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :
5.  $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$
6. **if**  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :
7. **return**  $Q_{\text{new}}$
8.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

"calculating"

## Q-LEARNING ( $\mathcal{S}, \mathcal{A}, \gamma, s_0, \alpha$ )

1. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :
2.  $Q_{\text{old}}(s, a) = 0$
3.  $s \leftarrow s_0$
4. **while** True:
5.  $a \leftarrow \text{select\_action}(s, Q_{\text{old}}(s, a))$
6.  $r, s' = \text{execute}(a)$
7.  $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$
8.  $s \leftarrow s'$
9. **if**  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :
10. **return**  $Q_{\text{new}}$
11.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

"estimating"

## Q-LEARNING ( $\mathcal{S}, \mathcal{A}, \gamma, s_0, \alpha$ )

1. **for**  $s \in \mathcal{S}, a \in \mathcal{A}$ :
2.  $Q_{\text{old}}(s, a) = 0$
3.  $s \leftarrow s_0$
4. **while** True:
5.  $a \leftarrow \text{select\_action}(s, Q_{\text{old}}(s, a))$
6.  $r, s' = \text{execute}(a)$
7.  $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$
8.  $s \leftarrow s'$
9. **if**  $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :
10. **return**  $Q_{\text{new}}$
11.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

- Remarkably, still can converge. (So long  $\mathcal{S}, \mathcal{A}$  are finite; we visit every state and action infinity-many times; and  $\alpha$  decays.)

- Line 7 :

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha \left( r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

is equivalently:

$$Q_{\text{new}}(s, a) \leftarrow Q_{\text{old}}(s, a) + \alpha \left( [r + \gamma \max_{a'} Q_{\text{old}}(s', a')] - Q_{\text{old}}(s, a) \right)$$

old belief + learning rate ( target - old belief )

pretty similar to SGD.

- Line 5, a sub-routine.



- If our  $Q$  values are estimated quite accurately (nearly converged to the true  $Q$  values), then should act greedily
  - $\arg \max_a Q^h(s, a)$ , as we did in MDP.
- During learning, especially in early stages, we'd like to explore.
- $\epsilon$ -greedy action selection strategy:
  - with probability  $1 - \epsilon$ , choose  $\arg \max_a Q(s, a)$
  - with probability  $\epsilon$ , choose an action  $a \in \mathcal{A}$  uniformly at random
- Benefit: get to observe more diverse  $(s,a)$  consequences.
- *exploration vs. exploitation.*

# Outline

- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?

- Q-learning only is kinda sensible for tabular setting.
- What do we do if  $\mathcal{S}$  and/or  $\mathcal{A}$  are large (or continuous)?
- Recall from Q-learning algorithm, key line 7 :

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha \left( r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

is equivalently:

$$Q_{\text{new}}(s, a) \leftarrow Q_{\text{old}}(s, a) + \alpha \left( [r + \gamma \max_{a'} Q_{\text{old}}(s', a')] - Q_{\text{old}}(s, a) \right)$$

old belief
+
learning rate
(
target
-
old belief
)

- Can be interpreted as we're minimizing:

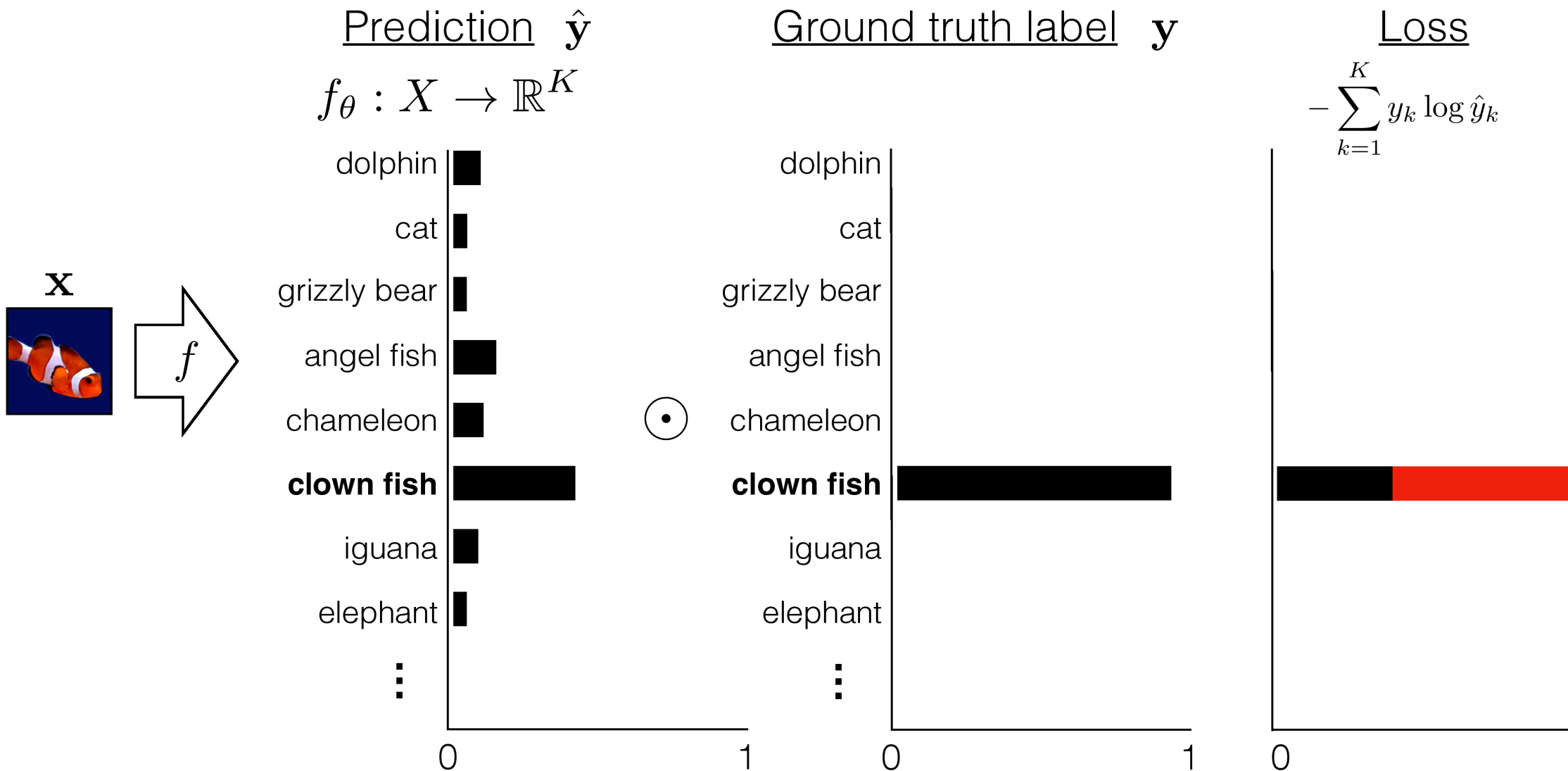
$$(Q(s, a) - (r + \gamma \max_{a'} Q(s', a')))^2$$

via gradient method!

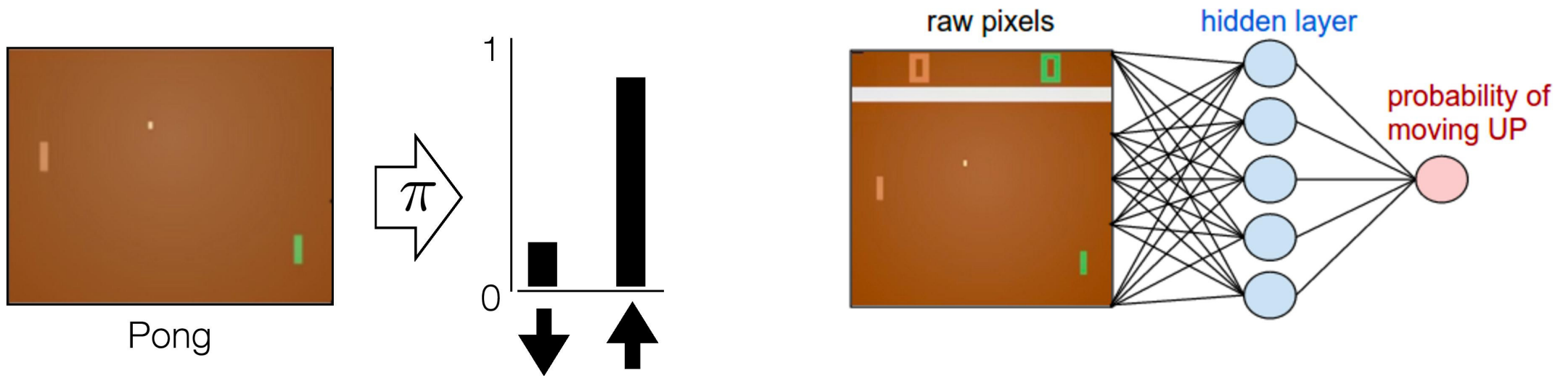
# Outline

- Recap: Markov Decision Processes
- Reinforcement Learning Setup
  - What's changed from MDP?
- Model-based methods
- Model-free methods
  - (tabular) Q-learning
    - $\epsilon$ -greedy action selection
    - exploration vs. exploitation
  - (neural network) Q-learning
- RL setup again
  - What's changed from supervised learning?

# Supervised learning

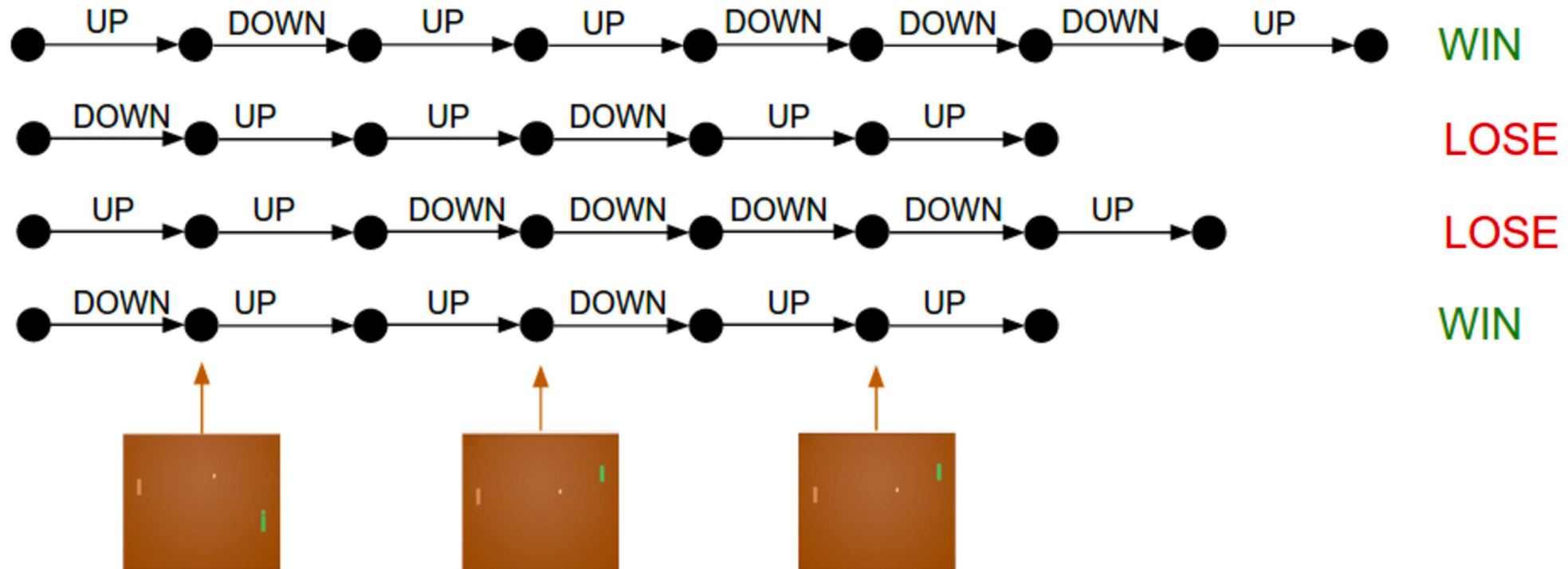


- If explicit “good” state-action pair is given, also supervised learning.
- Behavior cloning or imitation learning.
- But what if no explicit guide?



[Adapted from Andrej Karpathy: <http://karpathy.github.io/2016/05/31/r1/>]

- If no direct supervision is available?
- Strictly RL setting. Interact, observe and get data. Use rewards/value as "coy" supervision signal.



[Adapted from Andrej Karpathy: <http://karpathy.github.io/2016/05/31/rl/>]

We'd appreciate your [feedback](#) on the lecture.

Thanks!