

# COMS W4701: Artificial Intelligence

## Lecture 12: RL Control

Tony Dear, Ph.D.

Department of Computer Science  
School of Engineering and Applied Sciences

# Today

---

- State-action (Q) values
- Behavior and target policies
- Exploration,  $\epsilon$ -greedy policies
- SARSA
- Q-learning

# Solving Sequential Decision Problems

	Evaluate a fixed policy $\pi$ : Solve for $V^\pi$	Learn optimal value function $V^*$ or optimal policy $\pi^*$
Dynamic Programming (known model $T, R$ )	<ul style="list-style-type: none"><li>• Formulate and solve linear system of equations</li><li>• Iterative policy evaluation</li></ul>	<ul style="list-style-type: none"><li>• Value iteration for <math>V^*</math>, followed by policy extraction for <math>\pi^*</math></li></ul>
Reinforcement Learning (no model)	<ul style="list-style-type: none"><li>• First-visit Monte Carlo</li><li>• Constant-<math>\alpha</math> Monte Carlo</li><li>• TD(0)</li></ul>	???

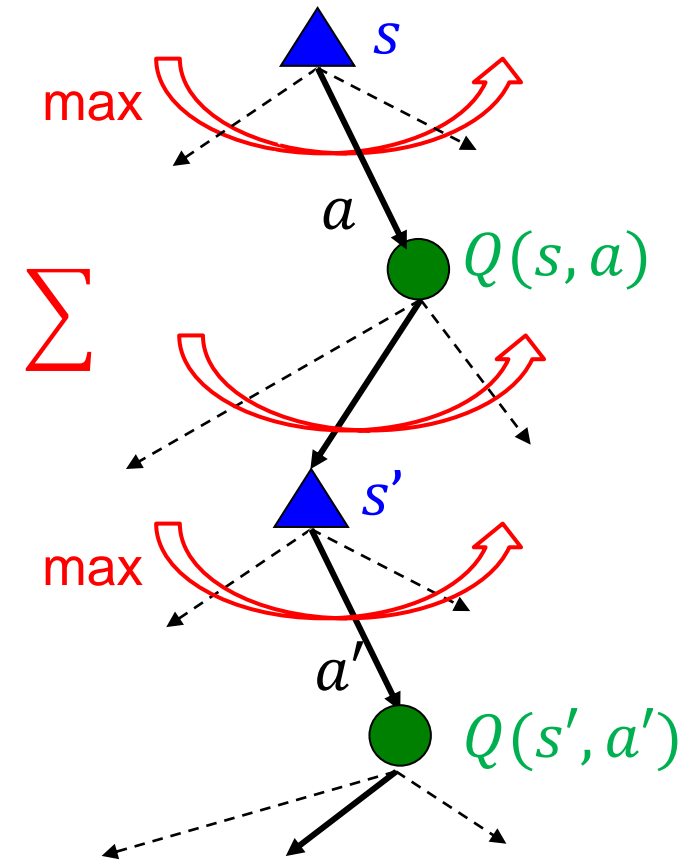
# State-Action Values

- How do we compute a policy  $\pi$  given its values  $V^\pi$ ?

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^\pi(s')]$$

don't have these!

- Idea: Directly *learn* **state-action values**  $Q^\pi(s, a)$  of the chance nodes using RL
- $Q^\pi(s, a)$  is the *expected* utility of  $s$  after taking action  $a$  and then following  $\pi$  thereafter

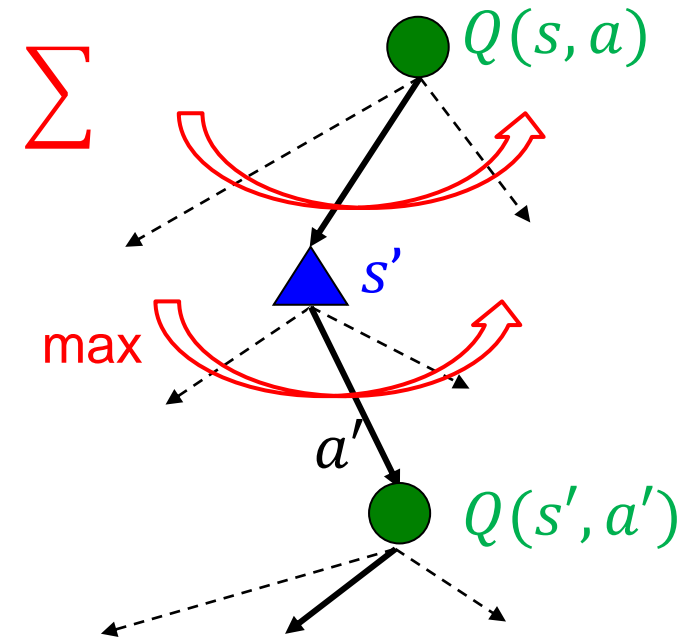


# Values and Policy

- **Q values** can be defined as functions of  $V^*$  or recursively as functions of themselves

$$\begin{aligned} Q(s, a) &= \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \\ &= \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q(s', a') \right] \end{aligned}$$

- To learn policies, we will first learn Q values
- Then use these to compute  $V^*$  or  $\pi^*$



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

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

# Example: Mini-Gridworld

- Suppose we have found the following  $Q$  values:

$A$	$B$	$C$
-----	-----	-----

$Q(s, a)$	$s = A$	$s = B$	$s = C$
$a = Left$	<b>4.061</b>	<b>4.364</b>	0.485
$a = Right$	1.152	2.364	<b>1.394</b>

- Value function:  $V^* = \max_a Q(s, a) = (4.061, 4.364, 1.394)$
- Policy:  $\pi^* = \operatorname{argmax}_a Q(s, a) = (Left, Left, Right)$

# Behavior Policies

---

- Suppose our agent is given some initial (possibly random) policy  $\pi$
- *Control* problem: We want to *improve* and make changes to  $\pi$  over time
- In RL, the agent always behaves according to the **behavior policy**
- E.g.,  $\pi(s) = \operatorname{argmax}_a Q(s, a)$  is a *greedy* behavior policy
- But this is only our **target policy** (optimal) when the Q-values are correct!
- How to behave if we are also learning Q-values at the same time?
- Instead of always acting greedily, add in some *exploration*

# $\epsilon$ -Greedy Policies

---

- **$\epsilon$ -greedy behavior policy:** Choose best action most of the time, but with small probability  $\epsilon$ , execute random action instead
- *Exploit* and choose action  $a = \operatorname{argmax}_{a'} Q(s, a')$  with probability  $1 - \epsilon$
- *Explore* and choose action uniformly at random with probability  $\epsilon$
- Exploration may lead to smaller *short-term* rewards, but crucial for discovering better actions and higher *long-term* rewards
- Exploration rate  $\epsilon$  is yet another tunable agent parameter



# TD Learning for Control

---

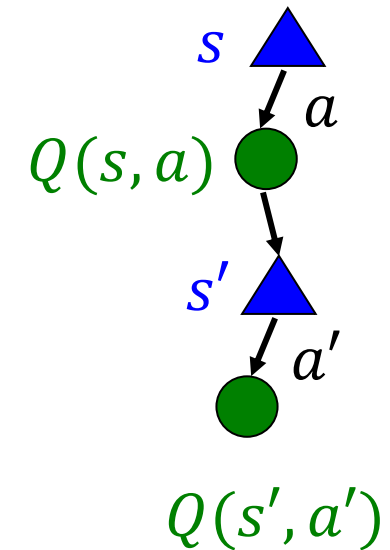
- $Q$  is just a more detailed version of  $V$ , so we can apply a TD approach to learn these values as well:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$$

- As in TD(0), we see the transition  $(s, a, r, s')$  and then update  $Q(s, a)$
- What do we use for  $a'$ , the action in the successor state  $s'$ ?
- In TD(0), the target uses  $V^\pi(s')$ , or the value of  $s'$  following  $\pi$
- So maybe we should select  $a'$  according to our *behavior policy* ( $\epsilon$ -greedy)
- Another possibility: Select  $a'$  according to our *target policy* (greedy)

# On-Policy Learning: SARSA

- **On-policy** learning: Q-value update is based on the successor Q-value corresponding to the *action that is actually taken according to the behavior policy  $\pi$*
- **Given:** Learning rate  $\alpha$ , exploration rate  $\varepsilon$ , discount factor  $\gamma$
- **Initialize**  $Q(s, a) \leftarrow 0$ , behavior policy  $\pi$  (e.g.,  $\varepsilon$ -greedy)
- **Loop:**
  - **Initialize** starting state  $s$ , action  $a = \pi(s)$  *if needed*
  - **Generate** sequence  $(s, a, r, s')$ ,  $a' \leftarrow \pi(s')$
  - $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$
  - $s \leftarrow s', a \leftarrow a'$



# Example: Mini-Gridworld

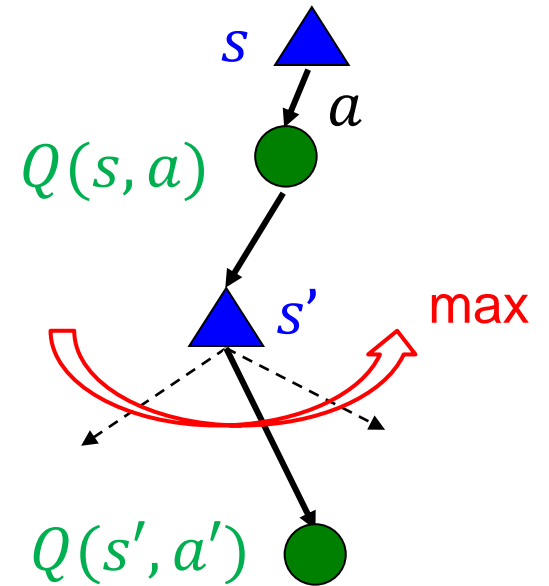
- Suppose we currently have  $Q(A, L) = 1$ ,  $Q(B, L) = 2$ ,  $Q(B, R) = 1.5$
- Behavior policy is  $\epsilon$ -greedy;  $\alpha = 0.5$ ,  $\gamma = 0.8$

$A$	$B$	$C$
-----	-----	-----

- Observed  $(s, a, r, s')$  sequence:  $A, L, +3, B$
- Suppose behavior policy generates  $a' = R$  (*explore*)
- Target:  $r + \gamma Q(B, R) = 3 + 0.8(1.5) = 4.2$
- Q-value update:  $Q(A, L) \leftarrow 1 + 0.5(4.2 - 1) = 2.6$

# Off-Policy Learning: Q-Learning

- **Off-policy** learning: Q-value update is based on the successor Q-value corresponding to the *action that “should be” taken according to the target policy*
- **Given:** Learning rate  $\alpha$ , exploration rate  $\varepsilon$ , discount factor  $\gamma$
- **Initialize**  $Q(s, a) \leftarrow 0$ , behavior policy  $\pi$  (e.g.,  $\varepsilon$ -greedy)
- **Loop:**
  - **Initialize** starting state  $s$  *if needed*, action  $a = \pi(s)$
  - **Generate** sequence  $(s, a, r, s')$
  - $Q(s, a) \leftarrow Q(s, a) + \alpha \left( r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)$
  - $s \leftarrow s'$



Action  $a'$  selected according to greedy policy ( $\varepsilon = 0$ )

# Example: Mini-Gridworld

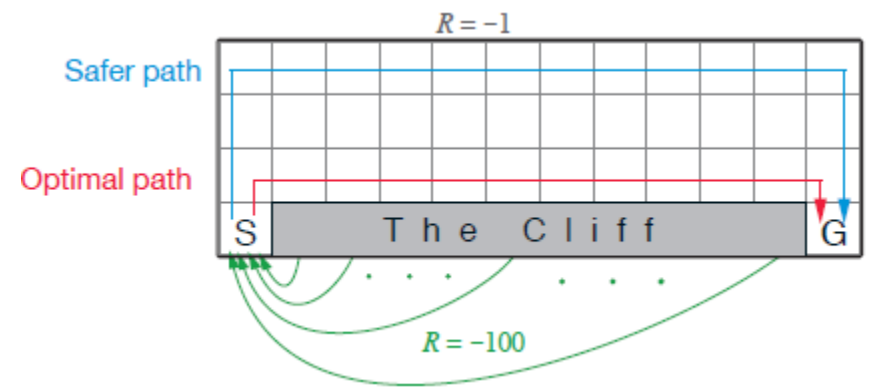
- Suppose we currently have  $Q(A, L) = 1$ ,  $Q(B, L) = 2$ ,  $Q(B, R) = 1.5$
- Behavior policy is  $\epsilon$ -greedy;  $\alpha = 0.5$ ,  $\gamma = 0.8$

$A$	$B$	$C$
-----	-----	-----

- Observed  $(s, a, r, s')$  sequence:  $A, L, +3, B$
- Target:  $r + \gamma \max_a Q(B, a) = r + \gamma Q(B, L) = 3 + 0.8(2) = 4.6$
- Q-value update:  $Q(A, L) \leftarrow 1 + 0.5(4.6 - 1) = 2.8$

# Example: Cliff Walking

- Deterministic grid world with one terminal (goal) state
- Living reward is  $-1$  in all states except for “cliff”, which rewards  $-100$
- Since transitions are deterministic, the optimal action at each state is to head in goal direction while ignoring the cliff
- We can use TD control to learn Q values and the corresponding policy



# Example: Cliff Walking

- With Q-learning, all learned Q-values will reflect the assumption that the agent will always act greedily (take action to move toward goal)
- Agent prefers “optimal path” as  $\max_a Q(s, a)$  is highest next to the cliff
- With SARSA, learned Q-values will be overall lower, esp near cliff
- Reflect all instances when agent chose to “explore” and jump off the cliff for — 100 reward
- Learned policy will be to take “safer” path!



# SARSA vs Q-Learning

---

- SARSA learns the optimal behavior policy, e.g.  $\epsilon$ -greedy
- Each Q-value update just requires knowledge of the next action
- Learned values are generally lower, reflecting suboptimal actions taken
- Q-learning learns the optimal target policy, e.g. greedy
- Each Q-value update requires a *max* over successor state actions
- Learned values are *not* affected by suboptimal actions from exploration
- SARSA & Q-learning are the same if behavior & target policies are as well



# Solving Sequential Decision Problems

	<b>Evaluate a fixed policy <math>\pi</math>: Solve for <math>V^\pi</math></b>	<b>Learn optimal value function <math>V^*</math> or optimal policy <math>\pi^*</math></b>
<b>Dynamic Programming (known model <math>T, R</math>)</b>	<ul style="list-style-type: none"><li>• Formulate and solve linear system of equations</li><li>• Iterative policy evaluation</li></ul>	<ul style="list-style-type: none"><li>• Value iteration for <math>V^*</math>, followed by policy extraction for <math>\pi^*</math></li></ul>
<b>Reinforcement Learning (no model)</b>	<ul style="list-style-type: none"><li>• First-visit Monte Carlo</li><li>• Constant-<math>\alpha</math> Monte Carlo</li><li>• TD(0)</li></ul>	<ul style="list-style-type: none"><li>• SARSA</li><li>• Q-learning followed by max / argmax operations</li></ul>

# Summary

---

- The RL control problem involves learning optimal policies and values from data
- We need to explore to try new states and actions; keep track of *state-action (Q) values*
- Simple approach is a  $\varepsilon$ -greedy behavior policy
- We can learn and update Q values using temporal difference learning
- SARSA (on-policy): Update Q value using Q value of action taken in successor state
- Q-learning (off-policy): Update current Q value using *best* Q value in successor state