

# Lecture 11: Model-Free Control

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

May 18 & 20, 2020

# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Model-Free Reinforcement Learning

- Last lecture:
  - ▶ Model-free prediction
  - ▶ Estimate the value function of an *unknown* MDP
- This lecture:
  - ▶ Model-free control
  - ▶ Optimize the value function of an *unknown* MDP

# Uses of Model-Free Control

Some example problems that can be modeled as MDPs

- Elevator
- Parallel Parking
- Ship Steering
- Bioreactor
- Helicopter
- Aeroplane Logistics
- Robocup Soccer
- Quake
- Portfolio management
- Protein Folding
- Robot walking
- Game of Go

For most of these problems, either:

- MDP model is unknown, but experience can be sampled
- MDP model is known, but is too big to use, except by samples

Model-free control can solve these problems

# On & Off-Policy Learning

- **Behavior-policy:** determines which action to take, from which we determine the next state to visit (also called “sampling policy”)
- **Target-policy:** determines the action that appears to be best (also called “learning policy”)
- Goal in reinforcement learning:
  - ▶ improve the target policy
  - ▶ while using a behavior policy to ensure that we visit states often enough (exploration schemes such as  $\epsilon$ -greedy)
- **On-policy learning:** when the learning policy and the sampling policy are the same
- **Off-policy learning:** when the learning policy and the sampling policy are different

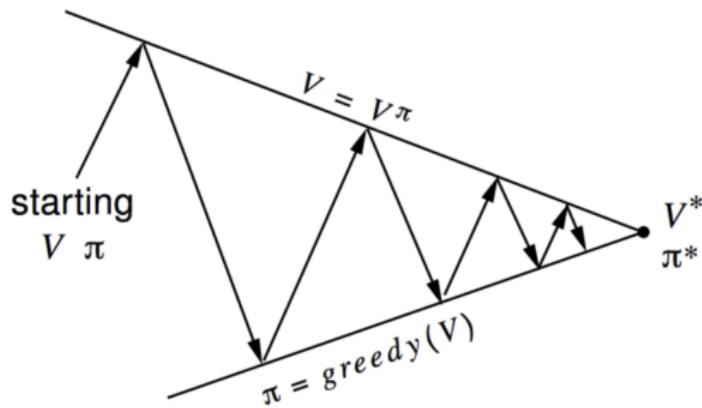
# On & Off-Policy Learning

- On-policy learning
  - ▶ "Learn on the job"
  - ▶ Learn the value of the target policy  $\pi$  from experience sampled from behavior policy  $\pi$
  - ▶ May not ensure enough exploration of state space
- Off-policy learning
  - ▶ "Look over someone's shoulder"
  - ▶ Learn the value of the target policy  $\pi$  from experience sampled from behavior policy  $\mu$
  - ▶ Learning is from experience(data) "off" the target policy
  - ▶ Compared to on-policy learning, off-policy learning is
    - ★ more powerful & general
    - ★ often of greater variance & slower to converge

# Outline

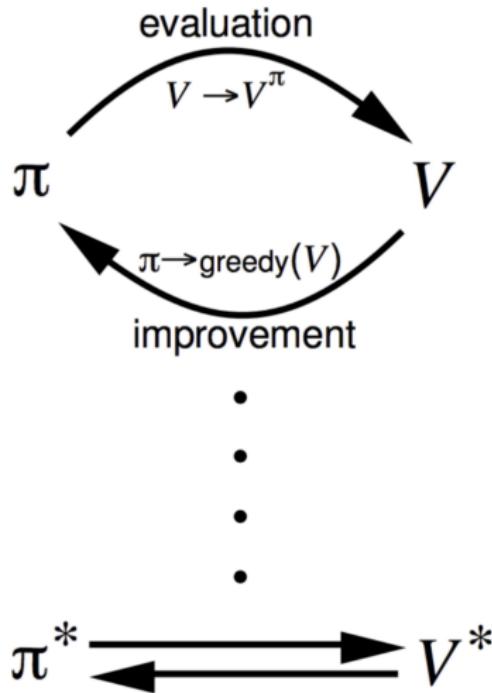
- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Generalized Policy Iteration (Refresher)

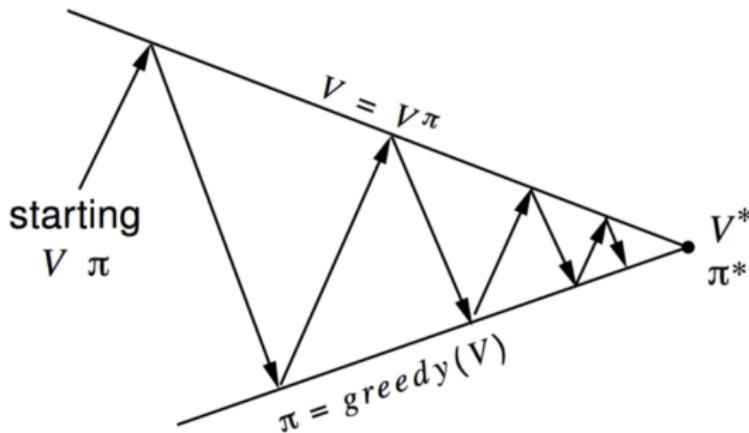


Policy evaluation Estimate  $v_\pi$   
e.g. Iterative policy evaluation

Policy improvement Generate  $\pi' \geq \pi$   
e.g. Greedy policy improvement



# Generalized Policy Iteration With Monte-Carlo Evaluation



Policy evaluation Monte-Carlo policy evaluation,  $V = v_\pi$ ?  
Policy improvement Greedy policy improvement?

# Model-Free Policy Iteration Using Action-Value Function

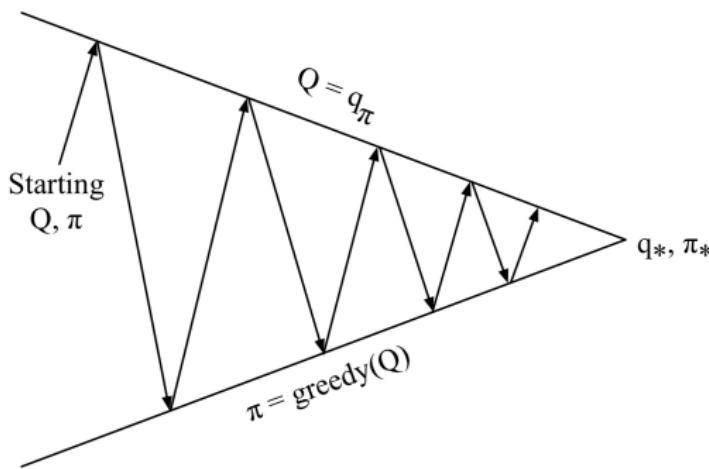
- Greedy policy improvement over  $V(s)$  requires model of MDP

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} \{ \mathcal{R}_s^a + \mathcal{P}_{ss'}^a V(s') \}$$

- Greedy policy improvement over  $Q(s, a)$  is model-free

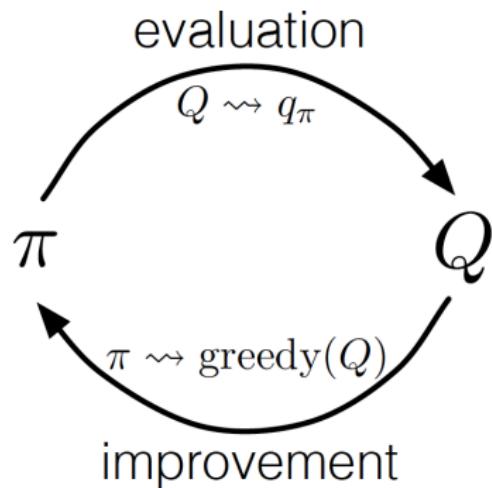
$$\pi'(s) = \arg \max_{a \in \mathcal{A}} Q(s, a)$$

# Generalized Policy Iteration with Action-Value Function



Policy evaluation Monte-Carlo policy evaluation,  $Q = q_\pi$   
Policy improvement Greedy policy improvement?

# Monte-Carlo Control



- **MC policy iteration:** policy evaluation using MC methods followed by policy improvement
- **Policy improvement:** greedify with respect to value (or action-value) function

# Example of Greedy Action Selection



"Behind one door is tenure - behind the other is flipping burgers at McDonald's."

- There are two doors in front of you.
- You open the left door and get reward 0  
 $V(\text{left}) = 0$
- You open the right door and get reward +1  
 $V(\text{right}) = +1$
- You open the right door and get reward +3  
 $V(\text{right}) = +3$
- You open the right door and get reward +2  
 $V(\text{right}) = +2$
- Are you sure you've chosen the best door?

# Monte Carlo ES (Exploring States)

Monte Carlo ES (Exploring Starts), for estimating  $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$  (arbitrarily), for all  $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose  $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $S_0, A_0$ , following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

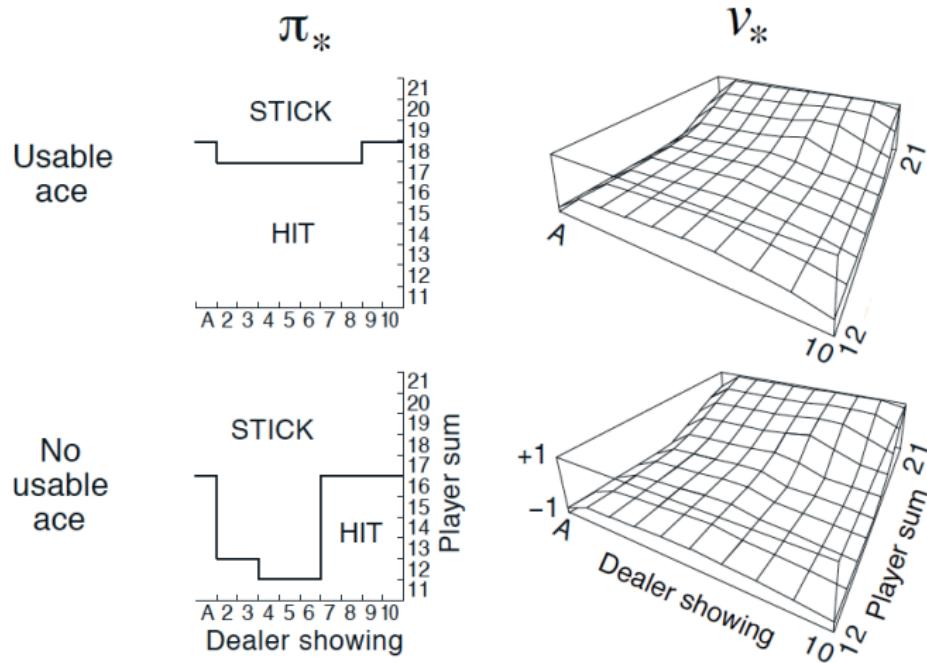
$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$

# Back to the Blackjack Example



# Monte-Carlo ES for Blackjack



# $\epsilon$ -Greedy Exploration

- How to avoid the unlikely assumption of exploring starts?
- Simplest idea for ensuring continual exploration
- All  $m$  actions are tried with non-zero probability
- With probability  $1 - \epsilon$  choose the greedy action
- With probability  $\epsilon$  choose an action at random

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a^* = \arg \max_{a \in \mathcal{A}} Q(s, a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

# $\epsilon$ -Greedy Policy Improvement

## Theorem

For any  $\epsilon$ -greedy policy  $\pi$ , the  $\epsilon$ -greedy policy  $\pi'$  with respect to  $q_\pi$  is an improvement,  $v_{\pi'}(s) \geq v_\pi(s)$

$$\begin{aligned}
 q_\pi(s, \pi'(s)) &= \underbrace{\sum_{a \in \mathcal{A}} \pi'(a|s) q_\pi(s, a)}_{\epsilon/m \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \max_{a \in \mathcal{A}} q_\pi(s, a)} \\
 &= \epsilon/m \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \max_{a \in \mathcal{A}} q_\pi(s, a) \\
 &\geq \epsilon/m \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \epsilon/m}{1 - \epsilon} q_\pi(s, a) \\
 &= \underbrace{\sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)}_{v_\pi(s)} = v_\pi(s)
 \end{aligned}$$

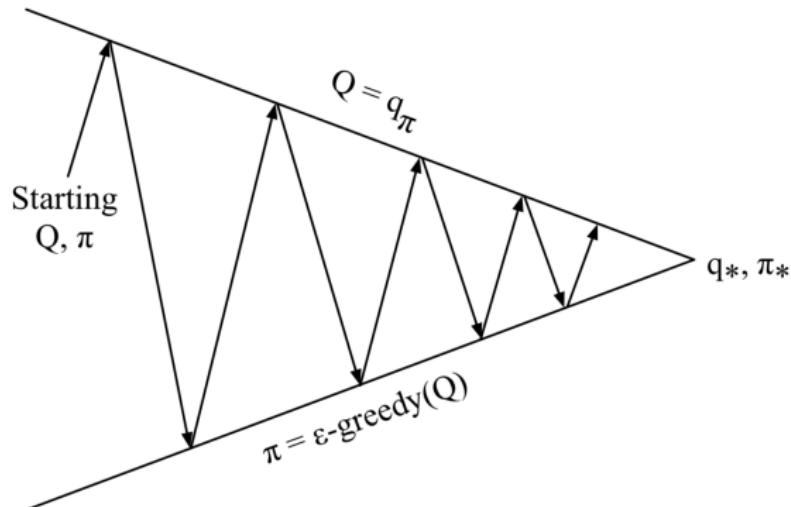
Therefore from policy improvement theorem,  $v_{\pi'}(s) \geq v_\pi(s)$

## Remark

$$1^{\circ} \cdot \underbrace{\sum_{a \in A} \frac{\pi(a|s) - \xi_m}{1-\varepsilon}}_{\text{underbrace}} = \frac{\sum_{a \in A} \pi(a|s) - \xi_{m,m}}{1-\varepsilon} = \frac{1-\varepsilon}{1-\varepsilon} = 1$$

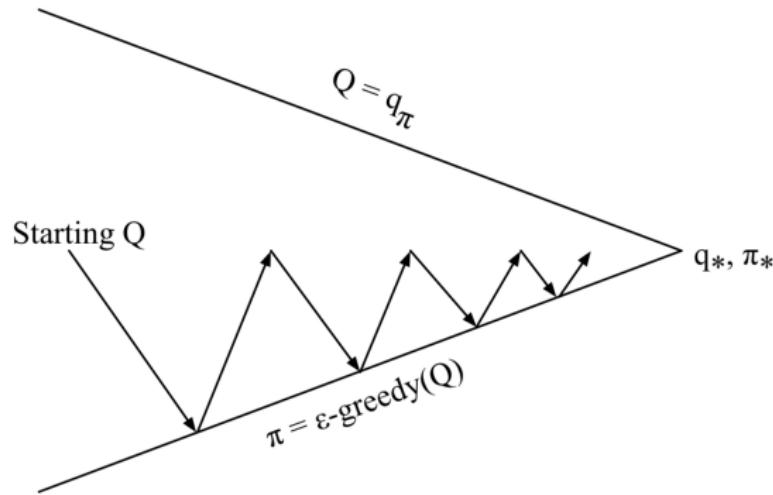
$$\begin{aligned} 2^{\circ} \cdot \underbrace{\sum_{a \in A} \frac{\pi(a|s) - \xi_m}{1-\varepsilon} \cdot q_\pi(s, a)}_{\text{underbrace}} &\leq \sum_{a \in A} \frac{\pi(a|s) - \xi_m}{1-\varepsilon} \cdot \underbrace{\max_{a \in A} q_\pi(s, a)}_{\text{underbrace}} \\ &= \left( \sum_{a \in A} \frac{\pi(a|s) - \xi_m}{1-\varepsilon} \right) \cdot \underbrace{\max_{a \in A} q_\pi(s, a)}_{\text{underbrace}} \\ &= \underbrace{\max_{a \in A} q_\pi(s, a)}_{\text{underbrace}} \end{aligned}$$

# Monte-Carlo Policy Iteration



Policy evaluation Monte-Carlo policy evaluation,  $Q = q_\pi$   
Policy improvement  $\epsilon$ -greedy policy improvement

# Monte-Carlo Control



Every episode:

Policy evaluation Monte-Carlo policy evaluation,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement

## Definition

Greedy in the Limit with Infinite Exploration (GLIE)

- All state-action pairs are explored infinitely many times,

$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty$$

- The policy converges on a greedy policy,

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbf{1}(a = \arg \max_{a' \in \mathcal{A}} Q_k(s, a'))$$

# GLIE

- For example,  $\epsilon$ -greedy is GLIE if  $\epsilon$  reduces to zero at  $\epsilon_k = \frac{1}{k}$
- Keep a declining exploration probability at a sufficiently slow rate that try all actions infinitely often
- In practice,  $\epsilon_k = \frac{1}{k^\beta}$  with  $\beta \in (0.5, 1]$

# GLIE Monte-Carlo Control

- Sample  $k$ th episode using  $\underline{\pi}$ :  $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$
- For each state  $S_t$  and action  $A_t$  in the episode,

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(G_t - Q(S_t, A_t))$$

- Improve policy based on new action-value function

$$\epsilon \leftarrow 1/k$$

$$\underline{\pi} \leftarrow \epsilon\text{-greedy}(Q)$$

## Theorem

*GLIE Monte-Carlo control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*

# On-Policy First-Visit MC Control

On-policy first-visit MC control (for  $\varepsilon$ -soft policies), estimates  $\pi \approx \pi_*$

Algorithm parameter: small  $\varepsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$  average( $Returns(S_t, A_t)$ )

$A^* \leftarrow \arg \max_a Q(S_t, a)$  (with ties broken arbitrarily)

For all  $a \in \mathcal{A}(S_t)$ :

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

# Partial Summary of MC

- MC has several advantages over DP
  - ▶ can learn directly from interaction with environment
  - ▶ No need for full models
  - ▶ No need to learn about ALL states (no bootstrapping)
  - ▶ Less harmed by violating Markov property
- MC methods provide an alternate policy evaluation process
- One issue to watch for: maintaining sufficient exploration
  - ▶ exploring starts, soft policies

# Boltzmann Exploration

- A limitation of  $\epsilon$ -greedy: when we explore, we choose actions at random without regard to their estimated values
- For larger action spaces, this can mean that we are spending a lot of time evaluating actions that are quite poor
- Boltzmann Exploration (also known as “Gibbs sampling” & “soft-max”): choosing an action based on its estimated value.

$$\pi(a|s) = \frac{e^{\beta \hat{Q}(s,a)}}{\sum_{a'} e^{\beta \hat{Q}(s,a')}}$$

$$a_s^* = \underset{a \in A}{\operatorname{argmax}} \hat{Q}(s,a)$$

$$\pi(a^*(s)) = \frac{e^{\beta \hat{Q}(s,a_s^*)}}{\sum_{a'} e^{\beta \hat{Q}(s,a')}}$$

- $\hat{Q}(s, a)$  is an estimate of the value of being in state s and taking action a.
- $\beta$  is a tunable parameter
  - ▶  $\beta = 0$  produces a pure exploration policy
  - ▶  $\beta \rightarrow \infty$  produces a greedy policy

$$= \frac{1 + \sum_{a \neq a^*} e^{\beta (\hat{Q}(s,a) - \hat{Q}(s,a^*))}}{1 + \sum_{a \neq a^*} e^{\beta (\hat{Q}(s,a) - \hat{Q}(s,a^*))}}$$

$\beta \rightarrow \infty$

# Outline

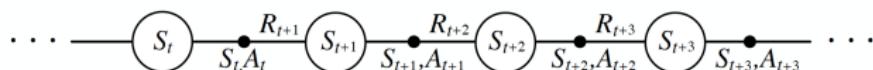
- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# MC vs. TD Control

- Temporal-difference (TD) learning has several advantages over Monte-Carlo (MC)
  - ▶ Lower variance
  - ▶ Online
  - ▶ Incomplete sequences
- Natural idea: use TD instead of MC in our control loop
  - ▶ Apply TD to  $Q(S, A)$
  - ▶ Use  $\epsilon$ -greedy policy improvement
  - ▶ Update every time-step

# Learning An Action-Value Function

Estimate  $q_\pi$  for the current policy  $\pi$



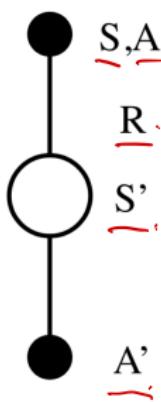
After every transition from a nonterminal state,  $S_t$ , do this:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \underline{Q(S_{t+1}, A_{t+1})} - Q(S_t, A_t)]$$

If  $S_{t+1}$  is terminal, then define  $Q(S_{t+1}, A_{t+1}) = 0$

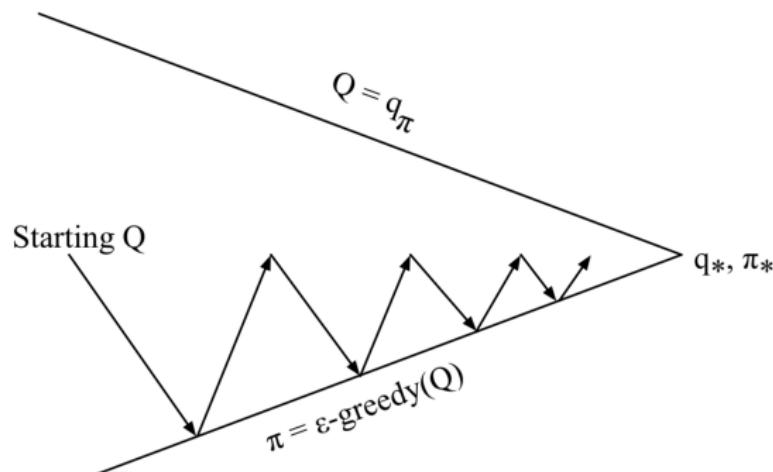
$$\overbrace{q_x(s, a) = \mathbb{E}_x [R_{t+1} + \gamma q_x(s_{t+1}, a_{t+1}) | S_t=s, A_t=a]}$$

# Updating Action-Value Functions with Sarsa



$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$$

# On-Policy Control With Sarsa



Every time-step:

Policy evaluation Sarsa,  $Q \approx q_\pi$

Policy improvement  $\epsilon$ -greedy policy improvement

# Sarsa Algorithm for On-Policy Control

Sarsa (on-policy TD control) for estimating  $Q \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

    Loop for each step of episode:

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S \leftarrow S'; A \leftarrow A'$ ;

    until  $S$  is terminal

# Convergence of Sarsa

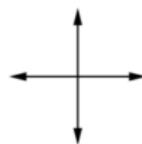
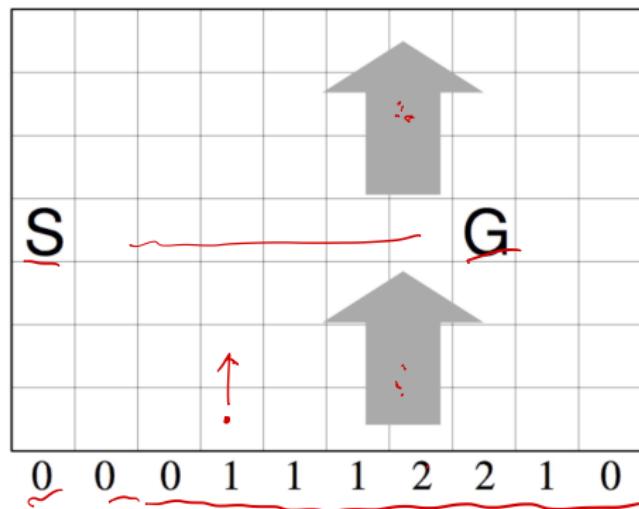
## Theorem

Sarsa converges to the optimal action-value function,  
 $Q(s, a) \rightarrow q_*(s, a)$ , under the following conditions:

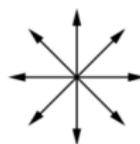
- GLIE sequence of policies  $\pi_t(a|s)$
- Robbins-Monro sequence of step-sizes  $\alpha_t$

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$
$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

# Windy Gridworld Example



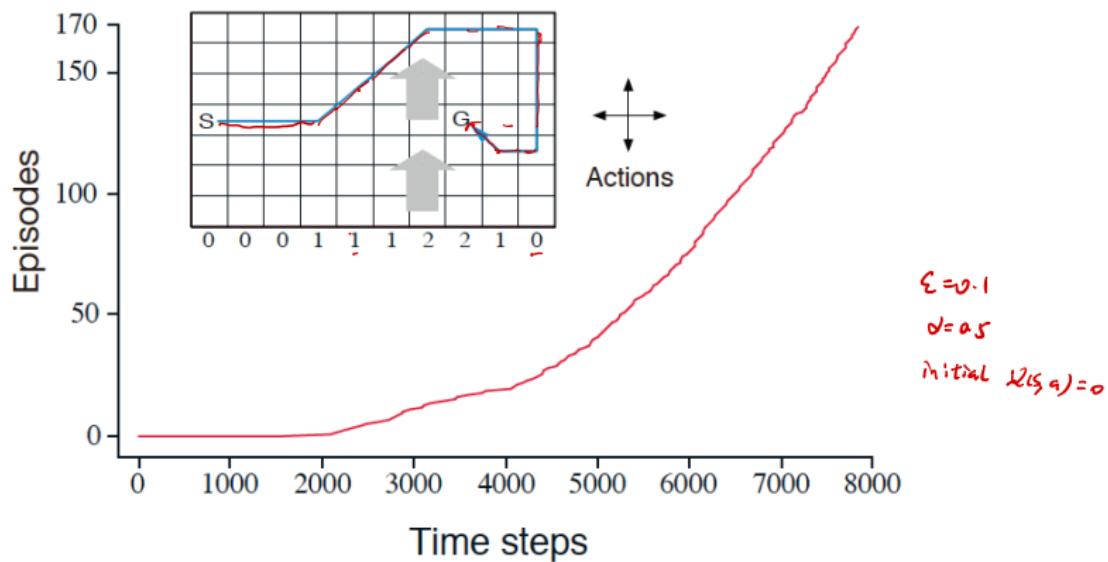
standard moves



king's moves

- Reward = -1 per time-step until reaching goal
- Undiscounted

# Sarsa on the Windy Gridworld



## *n*-Step Sarsa

- Consider the following *n*-step returns for  $n = 1, 2, \infty$ :

$$n = 1 \quad (\text{Sarsa}) \quad q_t^{(1)} = \underline{R_{t+1} + \gamma Q(S_{t+1})}$$

$$n = 2 \quad q_t^{(2)} = R_{t+1} + \gamma \underline{R_{t+2} + \gamma^2 Q(S_{t+2})}$$

⋮

⋮

$$n = \infty \quad (\text{MC}) \quad q_t^{(\infty)} = \underline{R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T}$$

- Define the *n*-step Q-return

$$\underline{q_t^{(n)}} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

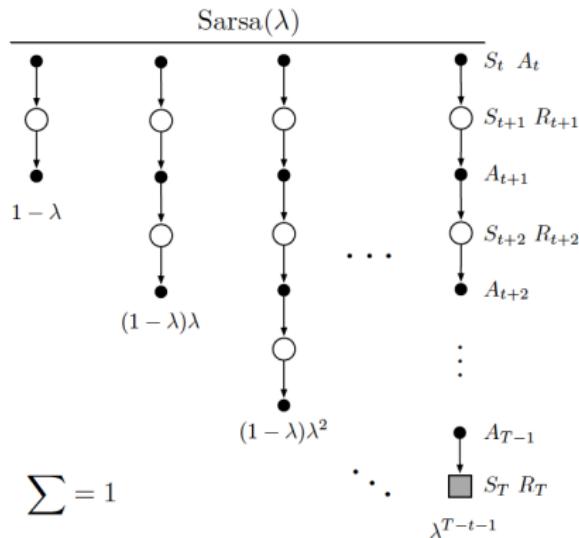
- n*-step Sarsa updates  $Q(s, a)$  towards the *n*-step Q-return

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( \underline{q_t^{(n)} - Q(S_t, A_t)} \right)$$

# Forward View Sarsa( $\lambda$ )

- Forward-view Sarsa( $\lambda$ )

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( \underline{q_t^\lambda} - Q(S_t, A_t) \right)$$



- The  $q^\lambda$  return combines all  $n$ -step Q>Returns  $q_t^{(n)}$
- Using weight  $(1 - \lambda)\lambda^{n-1}$

$$q_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$

## Backward View Sarsa( $\lambda$ )

- Just like TD( $\lambda$ ), we use eligibility traces in an online algorithm
- But Sarsa( $\lambda$ ) has one eligibility trace for each state-action pair

$$E_0(s, a) = 0$$

$$E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + \mathbf{1}(S_t = s, A_t = a)$$

- $Q(s, a)$  is updated for every state  $s$  and action  $a$
- In proportion to TD-error  $\delta_t$  and eligibility trace  $E_t(s, a)$

$$\delta_t = \underline{R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)}$$

$$Q(s, a) \leftarrow Q(s, a) + \underline{\alpha \delta_t E_t(s, a)}$$

# Sarsa( $\lambda$ ) Algorithm

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat (for each episode):

$$E(s, a) = 0, \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}(s)$$

Initialize  $S, A$

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$$\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$$

$$E(S, A) \leftarrow E(S, A) + 1$$

For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :

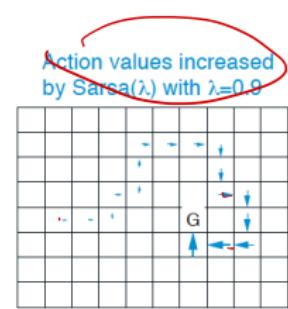
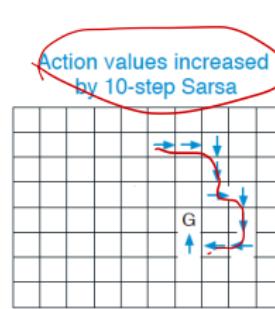
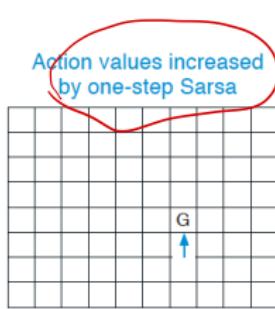
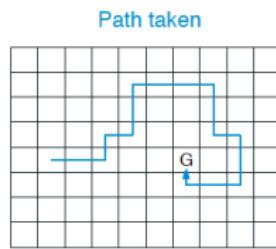
$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$$

$$E(s, a) \leftarrow \gamma \lambda E(s, a)$$

$$S \leftarrow S'; A \leftarrow A'$$

until  $S$  is terminal

# Sarsa( $\lambda$ ) Gridworld Example



# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Off-Policy Learning

- Learn the value of the target policy  $\pi$  from experience due to behavior policy  $\mu$
- In this sense, learning is from experience(data) “off” the target policy, and the overall process is termed off-policy learning.
- For example,  $\pi$  is the greedy policy (and ultimately the optimal policy) while  $\mu$  is exploratory (e.g.,  $\epsilon$ -soft)
- In general, we only require coverage, i.e., that  $\mu$  generates behavior that covers, or includes,  $\pi$

$$\mu(a|s) > 0 \text{ for every } s, a \text{ at which } \pi(a|s) > 0$$

# Off-Policy Learning

- Evaluate target policy  $\pi(a|s)$  to compute  $v_\pi(s)$  or  $q_\pi(s, a)$
- While following behaviour policy  $\mu(a|s)$

$$\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$$

- Why is this important?
  - ▶ Learn from observing humans or other agents
  - ▶ Re-use experience generated from old policies  $\pi_1, \pi_2, \dots, \pi_{t-1}$
  - ▶ Learn about *optimal* policy while following *exploratory* policy
  - ▶ Learn about *multiple* policies while following *one* policy

# Importance Sampling

- Estimate the expectation of a function

$$\mathbb{E}_{X \sim \pi}[g(X)] = \sum \pi(x) f(x) \approx \frac{1}{n} \sum_{k=1}^n g(x_k), x_k \sim \pi$$

- But sometimes it is difficult to sample  $x$  from  $\pi$ , then we can sample  $x$  from another distribution  $\mu$ , then correct the weight

$$\begin{aligned}\mathbb{E}_{X \sim \pi}[g(X)] &= \sum \pi(x) g(x) = \sum \mu(x) \frac{\pi(x)}{\mu(x)} g(x) \\ &= \mathbb{E}_{X \sim \mu} \left[ \frac{\pi(X)}{\mu(X)} g(X) \right] \approx \frac{1}{n} \sum_{k=1}^n \frac{\pi(x_k)}{\mu(x_k)} g(x_k), x_k \sim \mu\end{aligned}$$

- For off-policy learning: weight each return by the ratio of the probabilities of the trajectory under the two policies

# Importance Sampling for Off-Policy Monte-Carlo

- Use returns generated from  $\mu$  to evaluate  $\pi$
- Weight return  $G_t$  according to similarity between policies
- Multiply importance sampling corrections along whole episode

$$G_t^{\pi/\mu} = \left[ \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} \right] G_t$$

- Update value towards corrected return

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{\pi/\mu} - V(S_t) \right)$$

- Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
- Importance sampling can dramatically increase variance

# Importance Sampling Ratio

- Probability of the rest of the trajectory, after  $S_t$ , under  $\pi$

$$\begin{aligned} \Pr \{ A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi \} & \xrightarrow{\text{Chain rule + Markov property}} \\ & = \pi(A_t|S_t) p(S_{t+1}|S_t, A_t) \pi(A_{t+1}|S_{t+1}) \cdots p(S_T|S_{T-1}, A_{T-1}) \\ & = \prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{k+1}|S_k, A_k), \end{aligned}$$

- In importance sampling, each arm is weighted by the relative probability of the trajectory under the two policies

$$\rho_t^T = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} \mu(A_k|S_k) p(S_{k+1}|S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$$

*no MDP knowledge* p(S<sub>k+1</sub>|S<sub>k</sub>, A<sub>k</sub>)

- This is called the *importance sampling ratio*
- All importance sampling ratios have expected value 1

$$\mathbb{E}_{A_k \sim \mu} \left[ \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} \right] = \sum_a \mu(a|S_k) \frac{\pi(a|S_k)}{\mu(a|S_k)} = \sum_a \pi(a|S_k) = 1.$$

# Off-Policy MC Policy Evaluation

Incremental off-policy every-visit MC policy evaluation (returns  $Q \approx q_\pi$ )

Input: an arbitrary target policy  $\pi$

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

Repeat forever:

$\mu \leftarrow$  any policy with coverage of  $\pi$

Generate an episode using  $\mu$ :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For  $t = T - 1, T - 2, \dots$  downto 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$W \leftarrow W \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)}$$

If  $W = 0$  then ExitForLoop

# Off-Policy MC Control

Off-policy every-visit MC control (returns  $\pi \approx \pi_*$ )

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$\mu$   $\leftarrow$  any soft policy

Generate an episode using  $\mu$ :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For  $t = T - 1, T - 2, \dots$  downto 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If  $A_t \neq \pi(S_t)$  then ExitForLoop

$$W \leftarrow W \frac{1}{\mu(A_t | S_t)}$$

Target policy is greedy  
and deterministic

Behavior policy is soft,  
typically  $\epsilon$ -greedy

# Importance Sampling for Off-Policy TD

- Use TD targets generated from  $\mu$  to evaluate  $\pi$
- Weight TD target  $R + \gamma V(S')$  by importance sampling
- Only need a single importance sampling correction

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

- Much lower variance than Monte-Carlo importance sampling
- Policies only need to be similar over a single step

# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Q-Learning

- We now consider off-policy learning of action-values  $Q(s, a)$
- No importance sampling is required
- Next action to evaluate is chosen using behavior policy  
 $A_{t+1} \sim \mu(\cdot | S_t)$
- But we consider alternative successor action  $A' \sim \pi(\cdot | S_t)$
- Which means a separate policy is used to choose the alternative action in the future
- And update  $Q(S_t, A_t)$  towards value of alternative action

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \underbrace{Q(S_{t+1}, A') - Q(S_t, A_t)}_{})$$

# Off-Policy Control with Q-Learning

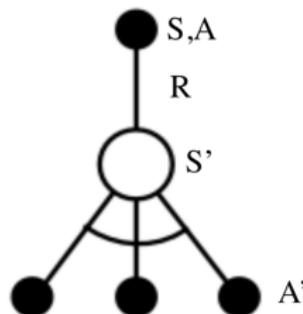
- We now allow both behaviour and target policies to improve
- The target policy  $\pi$  is greedy w.r.t.  $Q(s, a)$

$$\pi(S_{t+1}) = \arg \max_{a'} Q(S_{t+1}, a')$$

- The behaviour policy  $\mu$  is e.g.  $\epsilon$ -greedy w.r.t.  $Q(s, a)$
- The Q-learning target then simplifies:

$$\begin{aligned} & R_{t+1} + \gamma Q(S_{t+1}, A') \\ &= R_{t+1} + \gamma Q(S_{t+1}, \arg \max_{a'} Q(S_{t+1}, a')) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned}$$

# Q-Learning Control Algorithm



Bellman optimality Equation

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') \mid s_t=s, a_t=a]$$

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

## Theorem

*Q-learning control converges to the optimal action-value function,*  
 $Q(s, a) \rightarrow q_*(s, a)$

# Q-Learning Algorithm for Off-Policy Control

Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$S \leftarrow S'$

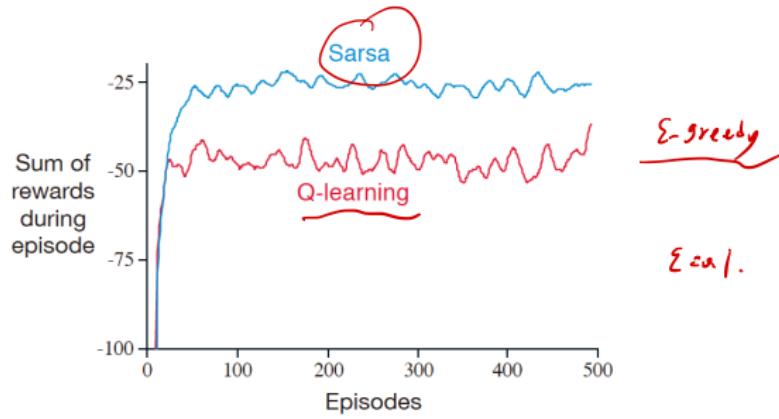
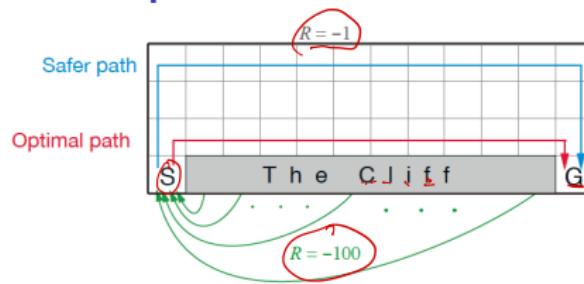
    until  $S$  is terminal

# Q-Learning Demo

<https://www.cs.ubc.ca/~poole/demos/rl/q.html>



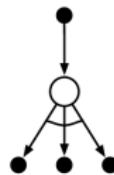
# Cliff Walking Example



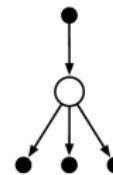
# Expected Sarsa

- Instead of the sample value-of-next-state, use the expectation!
- Expected Sarsa's performs better than Sarsa (but costs more)

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right] \end{aligned}$$



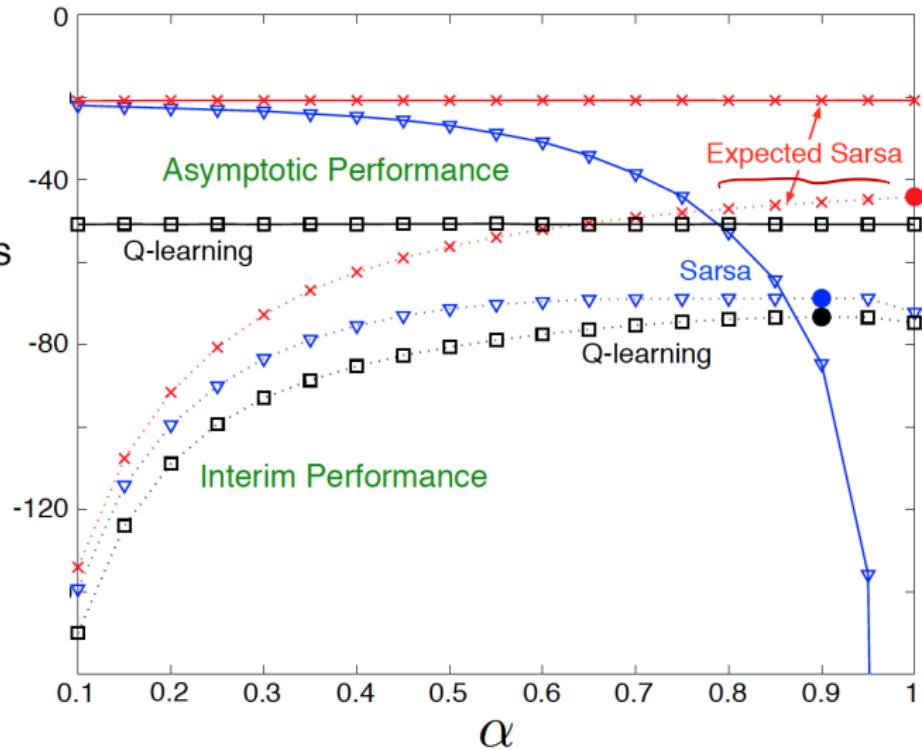
Q-learning



Expected Sarsa

# Performance on Cliff Walking Example

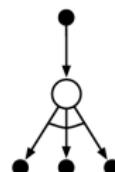
Sum of rewards  
per episode



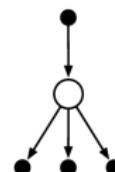
# Off-Policy Expected Sarsa

- Expected Sarsa generalizes to arbitrary behavior policies  $\mu$ 
  - ▶ in which case it includes Q-learning as the special case in which  $\pi$  is the greedy policy
- This idea seems to be new

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \underbrace{\pi(a|S_{t+1})Q(S_{t+1}, a)}_{\text{x-greedy}} - Q(S_t, A_t) \right] \end{aligned}$$



Q-learning



Expected Sarsa

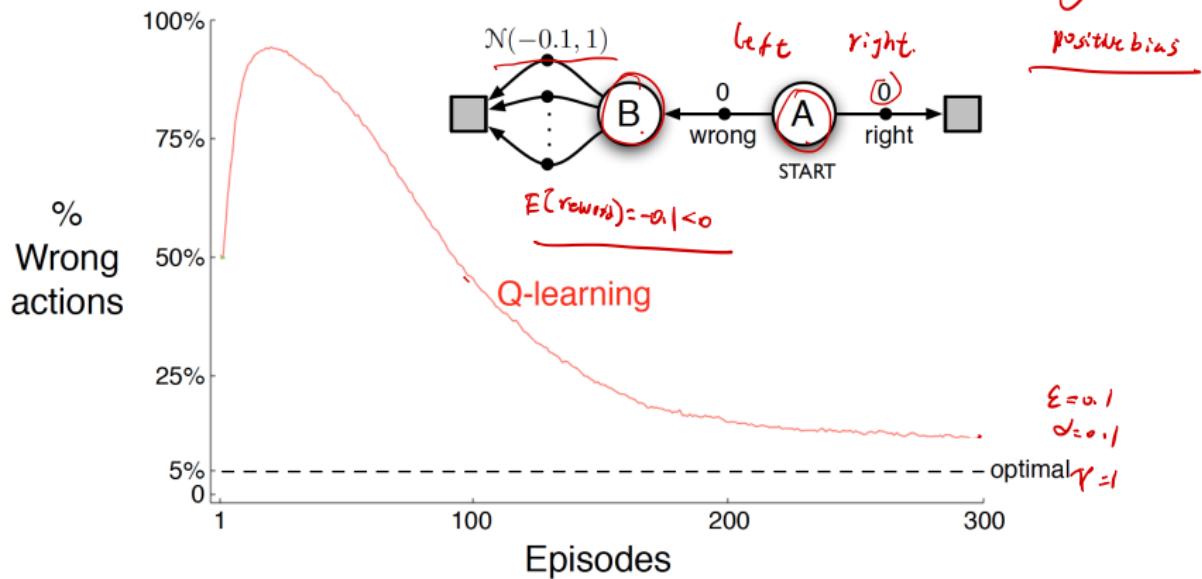
$\max_{a'} Q(S_{t+1}, a')$

# Maximization Bias Example

Maximization of estimation



positive bias



**Tabular Q-learning:** 
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

# Solution: Double Q-Learning

- Train 2 action-value functions  $Q_1$  &  $Q_2$
- Do Q-learning on both, but
  - ▶ never on the same time steps ( $Q_1$  &  $Q_2$  are independent)
  - ▶ pick  $Q_1$  or  $Q_2$  at random to be updated on each step

- If updating  $Q_1$ , use  $Q_2$  for the value of the next state:

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left( R_{t+1} + Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right)$$

- Action selections are (say)  $\epsilon$ -greedy with respect to the sum of  $Q_1$  and  $Q_2$

Applied to DQN Variants

# Solution: Double Q-Learning

Double Q-learning, for estimating  $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , such that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using the policy  $\varepsilon$ -greedy in  $Q_1 + Q_2$

        Take action  $A$ , observe  $R, S'$

        With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left( R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

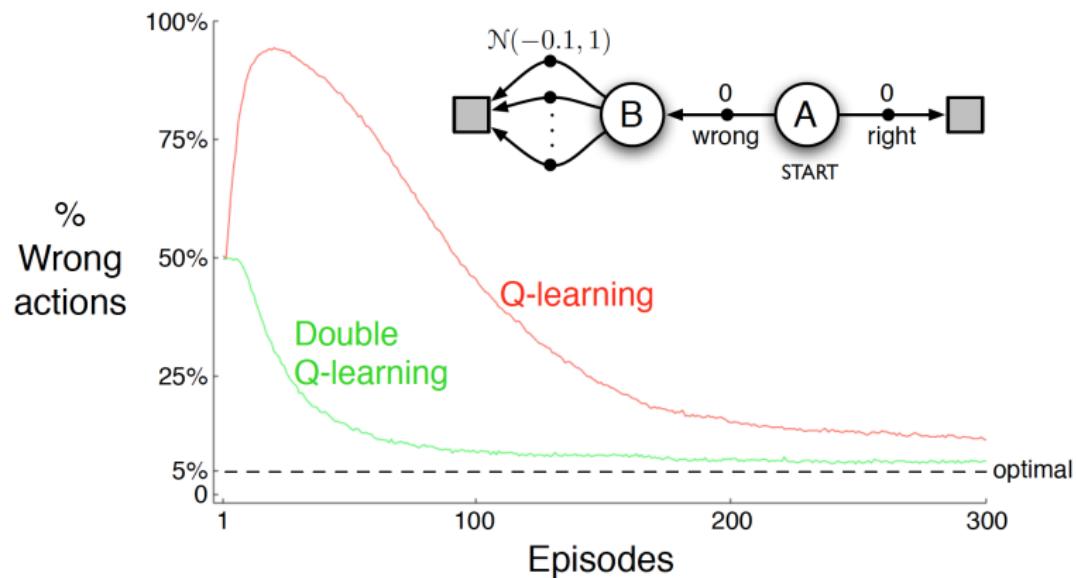
        else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left( R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$

    until  $S$  is terminal

# Maximization Bias Example

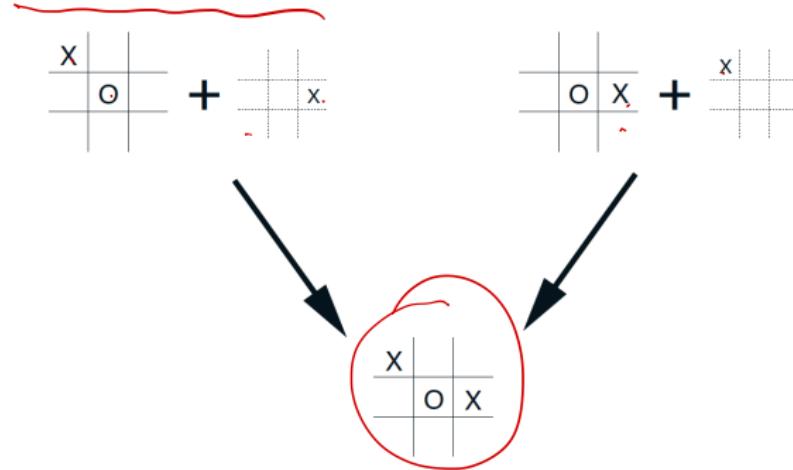


Double Q-learning:

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left( R_{t+1} + Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right)$$

# Afterstates

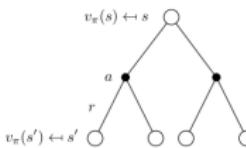
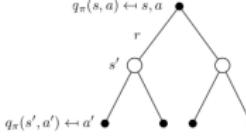
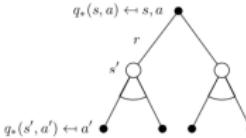
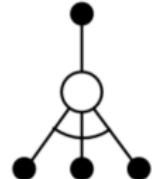
- Usually, a state-value function evaluates states in which the agent can take an action
- But sometimes it is useful to evaluate states after agent has acted, as in tic-tac-toe
- More efficient to design algorithms



# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

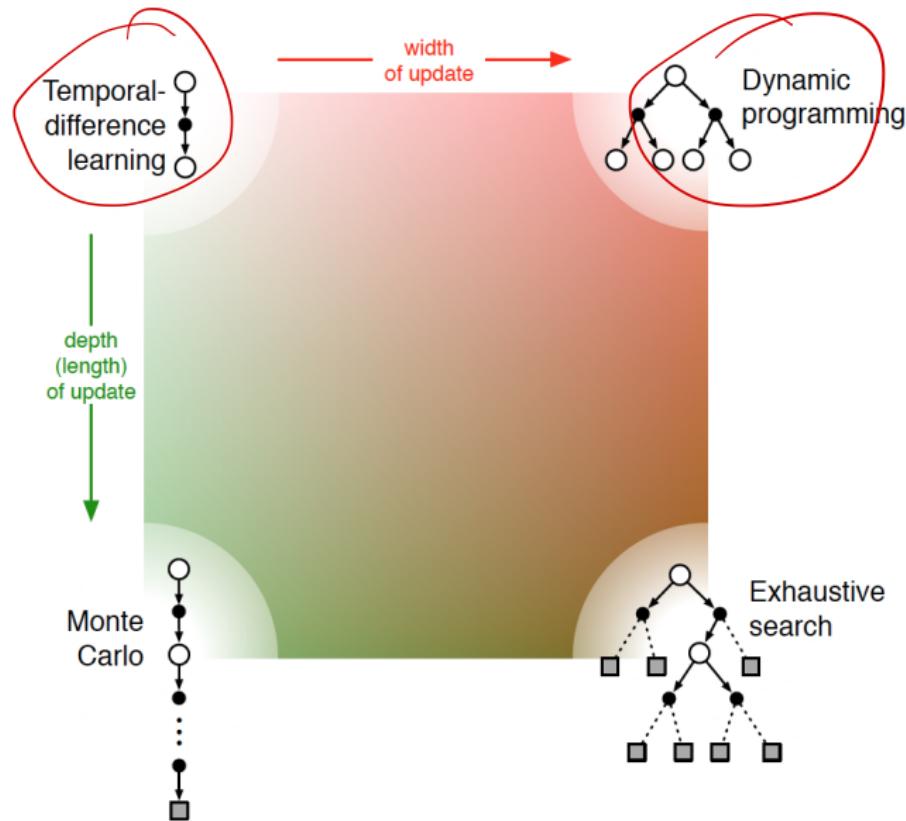
# Relationship Between DP and TD

	Full Backup (DP)	Sample Backup (TD)
Bellman Expectation Equation for $v_\pi(s)$	 <p><math>v_\pi(s) \leftarrow s</math></p> <p><math>v_\pi(s') \leftarrow s'</math></p> <p>Iterative Policy Evaluation</p>	
Bellman Expectation Equation for $q_\pi(s, a)$	 <p><math>q_\pi(s, a) \leftarrow s, a</math></p> <p><math>q_\pi(s', a') \leftarrow a'</math></p> <p>Q-Policy Iteration</p>	 <p>S, A</p> <p>R</p> <p>S'</p> <p>A'</p> <p>Sarsa</p>
Bellman Optimality Equation for $q_*(s, a)$	 <p><math>q_*(s, a) \leftarrow s, a</math></p> <p><math>q_*(s', a') \leftarrow a'</math></p> <p>Q-Value Iteration</p>	 <p>Q-Learning</p>

# Relationship Between DP and TD(2)

Full Backup (DP)		Sample Backup (TD)
Iterative Policy Evaluation $V(s) \leftarrow \mathbb{E}[R + \gamma V(S')   s]$		TD Learning $V(S) \xleftarrow{\alpha} R + \gamma V(S')$
Q-Policy Iteration $Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A')   s, a]$		Sarsa $Q(S, A) \xleftarrow{\alpha} R + \gamma Q(S', A')$
Q-Value Iteration $Q(s, a) \leftarrow \mathbb{E} \left[ R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')   s, a \right]$		Q-Learning $Q(S, A) \xleftarrow{\alpha} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')$

# Unified View of Reinforcement Learning



# Outline

- 1 Introduction
- 2 On-Policy Monte-Carlo Control
- 3 On-Policy Temporal-Difference Learning
- 4 Off-Policy Learning: Importance Sampling
- 5 Off-policy Learning: Q-learning
- 6 Summary
- 7 References

# Main References

- Reinforcement Learning: An Introduction (second edition), R. Sutton & A. Barto, 2018.
- RL course slides from Richard Sutton, University of Alberta.
- RL course slides from David Silver, University College London.