

# Reinforcement Learning

## Lecture 6

Model-free Control

Monte Carlo

TD: Sarsa and Q-Learning

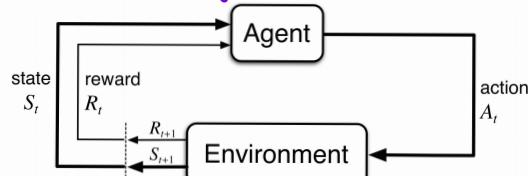
Robert Peharz

Institute of Theoretical Computer Science  
Graz University of Technology

Winter Term 2023/24

# Recap

policy  $\pi(a|s)$



MDP  $(\mathcal{S}, \mathcal{A}, p)$

greedy policy

$$\pi'(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q^{\pi}(s, a') \\ 0 & \text{otherwise} \end{cases}$$

**Policy Improvement Theorem**

$$V_{\pi}(s) \leq V_{\pi'}(s) \quad \forall s \in \mathcal{S}$$

**Theorem: Global Optimal Policy**

For any initial policy  $\pi_0$ , policy iteration converges to an optimal policy with value function  $v^*$

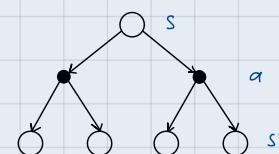
(Generalized) Policy Iteration

discounted return  $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma G_{t+1}$   
 $0 \leq \gamma \leq 1$

value function  $V_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s]$

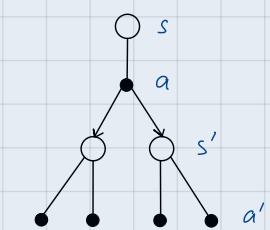
q-function  $q^{\pi}(s, a) := \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$

Bellman Expectation Equation in  $V$

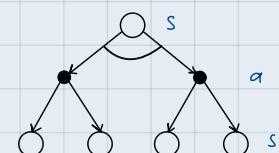


Bellman Expectation Equation in  $q$

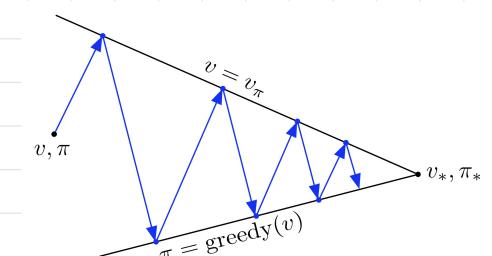
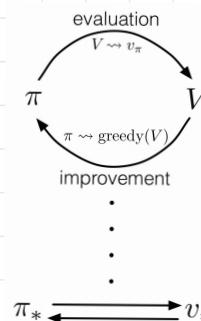
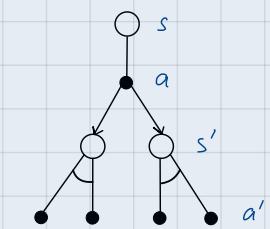
$$\begin{aligned} \hat{\Sigma} &\triangleq \sum \\ \hat{\max} &\triangleq \max \end{aligned}$$



Bellman Optimality Equation in  $V$



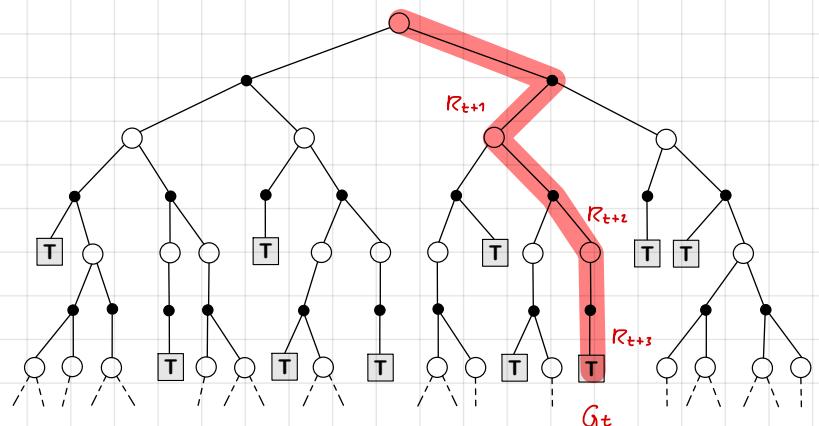
Bellman Optimality Equation in  $q$



# Recap: Model-free Prediction

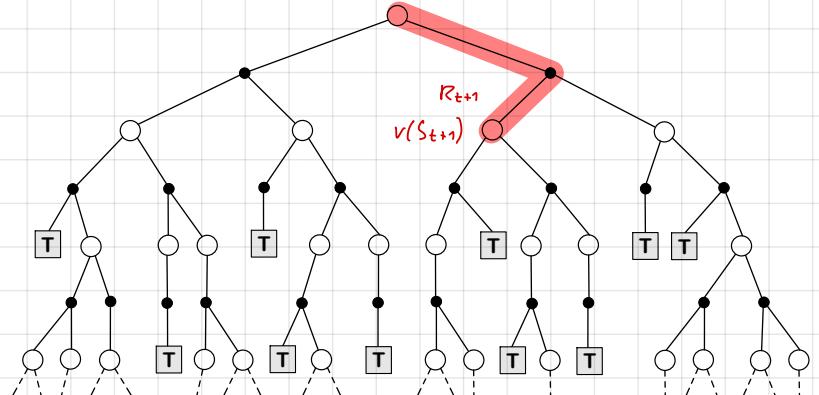
## Monte Carlo

$$v(s) \leftarrow v(s) + \frac{1}{N(s)} (G_t - v(s))$$

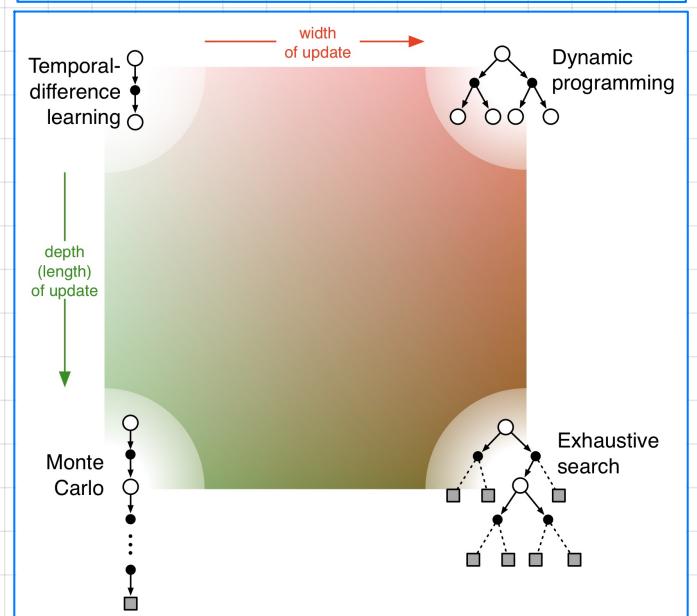
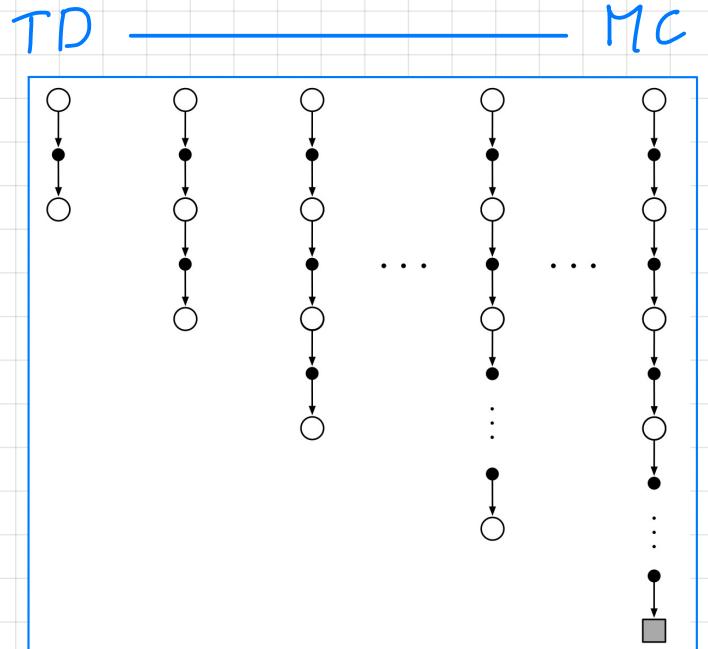


## Temporal Differences (TD)

$$v(s) \leftarrow v(s) + \alpha (R_{t+1} + \gamma v(s') - v(s))$$



## N-step TD



# **Model-free Control**

# Why Model-Free Control?

Images: Wikipedia

- MDP might be unknown

robotics

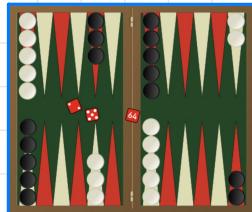


autonomous navigation



- We are too lazy to describe the MDP explicitly

games



- MDP too large

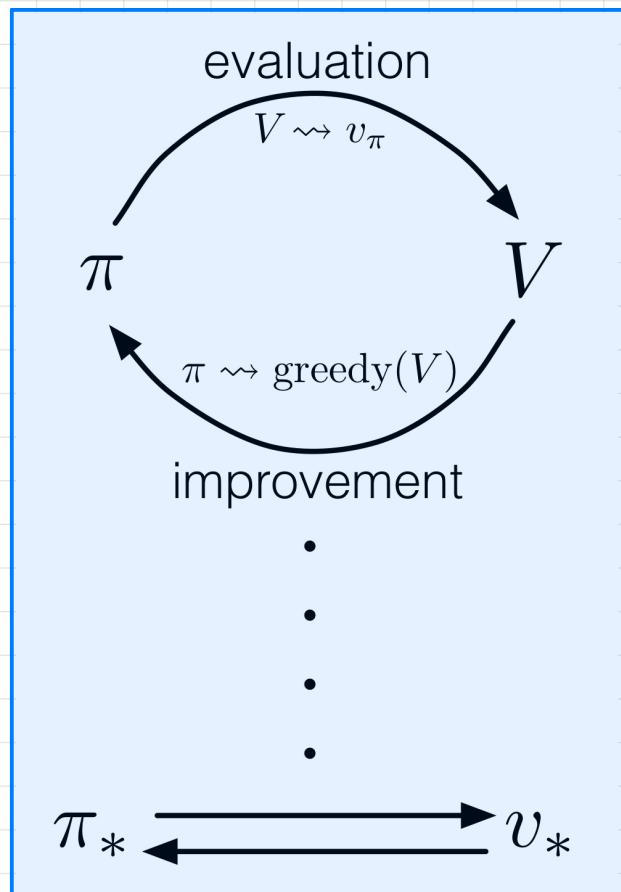
atari games



160 x 192 pixels  
128 colors

30720 states  
128

# Generalised Policy Iteration (GPI)



- Alternating between (approximately) computing  $v_{\tilde{\pi}}$  and (approximately) computing greedy( $v_{\tilde{\pi}}$ )
- Worked for model-based RL
- We know already two model-free techniques to compute  $v_{\tilde{\pi}}$ , MC and TD
- Essentially, we will plug them into GPI
- There are two subtleties, however

## Subtlety 1: We need q

- Policy improvement:

$$\hat{\pi}'(a|s) = \begin{cases} 1 & \text{if } a = \underset{a'}{\operatorname{argmax}} \ q_{\pi}(s, a') \\ 0 & \text{o.w.} \end{cases}$$

For planning, it was enough to estimate  $V_{\pi}(s)$ , since

$$q_{\pi}(s, a) = \underbrace{r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_{\pi}(s')}_{\text{MDP, now unknown}}$$

- Thus, we need to directly estimate  $q_{\pi}$
- No problem, MC and TD carry over to (s,a) pairs

# Basically, the same updates

- Monte Carlo

$$v(s) \leftarrow v(s) + \frac{1}{N(s)} (G_t - v(s))$$

$$q(s, a) \leftarrow q(s, a) + \frac{1}{N(s)} (G_t - q(s, a))$$

- Temporal Differences (TD)

$$v(s) \leftarrow v(s) + \alpha (R_{t+1} + \gamma v(s') - v(s))$$

$$q(s, a) \leftarrow q(s, a) + \alpha (R_{t+1} + \gamma q(s', a') - q(s, a))$$

- But, more expensive

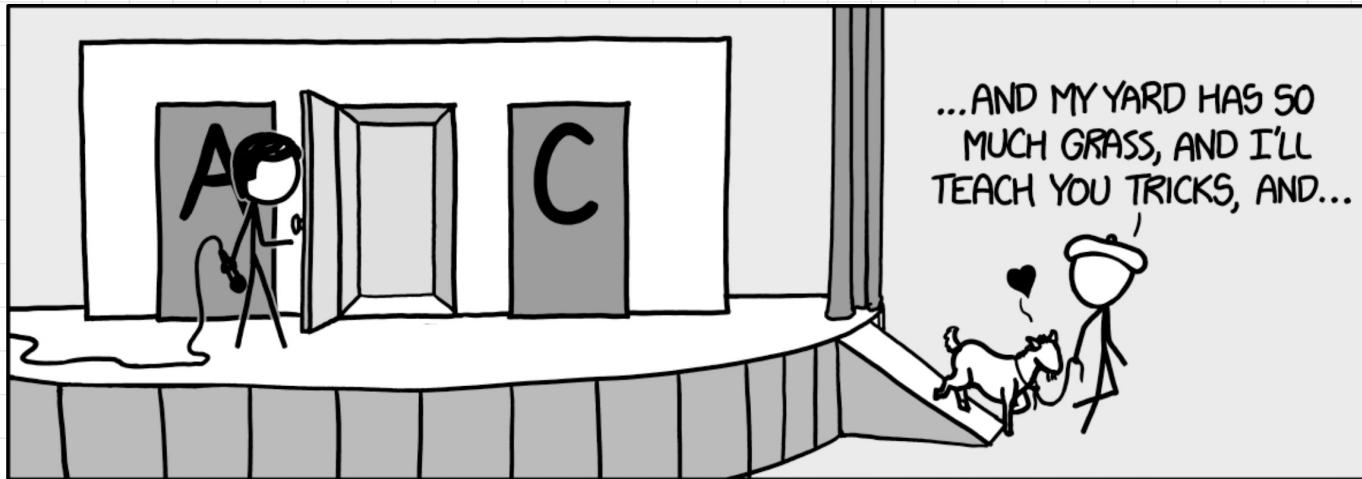
$v_{\pi} \dots |S'|$  values

$q_{\pi} \dots |S'| \times |A|$  values

$$\begin{pmatrix} v_{\pi}(s_1) \\ v_{\pi}(s_2) \\ \vdots \\ v_{\pi}(s_N) \end{pmatrix} \quad \begin{pmatrix} q_{\pi}(s_1, a_1) & \dots & q_{\pi}(s_1, a_M) \\ \vdots & \ddots & \vdots \\ q_{\pi}(s_N, a_1) & \dots & q_{\pi}(s_N, a_M) \end{pmatrix}$$

## Subtlety 2: Exploration

Image: xkcd.com



- You have the choice between doors A, B, and C.
  - you pick door A and get 5 reward
  - you pick door B and get 0 reward
  - you pick door C and get 1 reward
  - door A → 3 reward
  - door A → 4 reward
  - door A → 7 reward
- Is door A the best? → Exploration vs. Exploitation  
Fundamental trade-off in RL

# Monte Carlo Control with Exploring Starts

- initialize  $q(s, a)$  arbitrarily  $\forall s \in S, a \in A$
- $N(s, a) \leftarrow 0 \quad \forall s \in S, a \in A$  // visit counter
- repeat
  - select  $s_0, a_0$  randomly, with  $p(s, a) > 0, \forall s \in S, a \in A$  // assumes episodic tasks
  - generate episode  $s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T, s_T$  where  $a_t = \underset{a'}{\operatorname{argmax}} q(s_t, a')$  //  $\pi = \text{greedy}(q)$
  - $G \leftarrow 0$
  - for  $t = T-1, \dots, 0$ :
    - $G \leftarrow r_{t+1} + \gamma G$
    - $N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$
    - $q(s_t, a_t) \leftarrow q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G - q(s_t, a_t))$  // incremental average

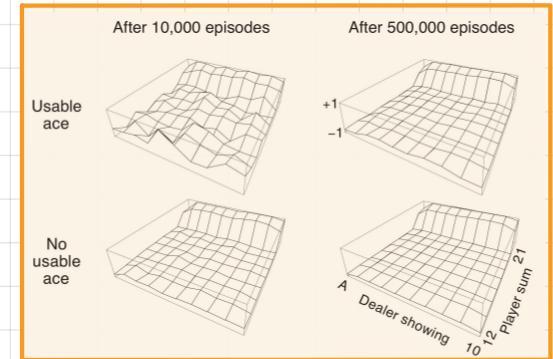
## Monte Carlo Control with Exploring Starts cont'd

- $N(s,a) \rightarrow \infty$  (since  $p(s,a) > 0$ )  $\forall s \in S, a \in A$
- Thus, eventually we will have fully explored the environment
- $\hat{\pi}$  cannot converge to suboptimal policy, since  $q$  would converge to  $q_{\hat{\pi}}$ , leading to a change of  $\hat{\pi}$ .
- Thus, if algorithm converges, it converges to  $\hat{\pi}^*$ ,  $q^*$  ✓
- Does it always converge? 
- Empirically, yes.
- Theoretically, open problem.
- Convergence with probability 1 has been proven for synchronous version, i.e. when we simulate an episode for each  $s,a$  pair, before updating  $q$  and  $\gamma < 1$ .

[Tsitsiklis, JMLR, 2002]

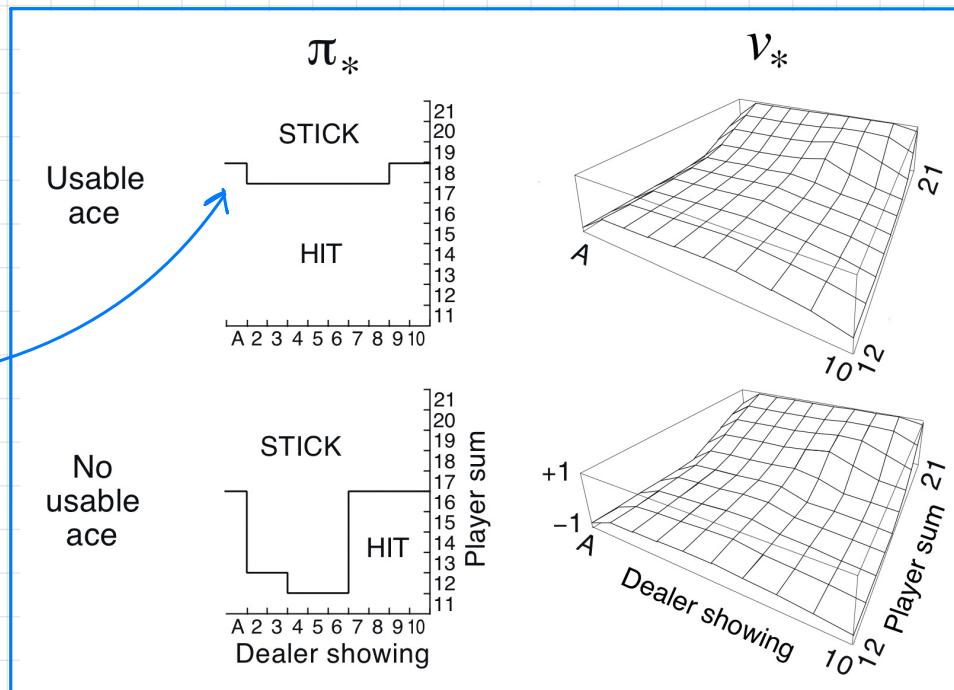
## Black Jack cont'd

Recall: Last lecture we evaluated Black Jack with (sub-optimal) "Hit-until-20" strategy



Using MC control, we can find an optimal policy:

"Basic strategy" by Thorp (1966), except for this notch



# Epsilon Greedy Policies

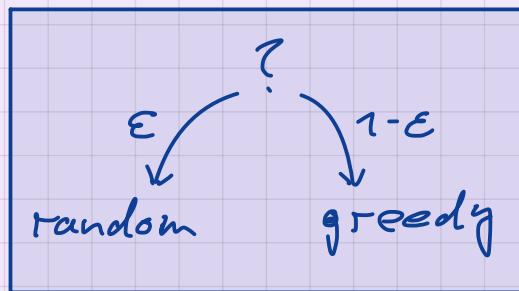
- Exploring starts are unrealistic / impractical...
- Instead, we can incorporate exploration in  $\pi$

$\epsilon$ -soft Policy:  $\hat{\pi}(a|s) \geq \epsilon$   $\forall s, a$



Image: [pinterest.uk](#)

$\epsilon$ -greedy Policy: With probability  $\epsilon$ , take random action (uniformly), otherwise take greedy action



$$\hat{\pi}(a|s) = \begin{cases} 1-\epsilon + \frac{\epsilon}{|\mathcal{A}(s)|} & \text{if } a = \arg\max_{a'} q(s, a') \\ \frac{\epsilon}{|\mathcal{A}(s)|} & \text{o.w.} \end{cases}$$

- $\epsilon$ -greedy: the "hardest"  $\epsilon$ -soft policy
- for  $\epsilon=0$ : greedy policy

## Greedy in the Limit with Infinite Exploration

- All state-action pairs are visited infinitely often
- Policy converges to greedy policy

In the GLIE setting, if MC control converges,\* it converges to  $\pi^*$  and  $q^*$

\* it probably always does

- assume that any state  $s'$  can be reached from any state  $s$
- then  $\epsilon$ -greedy policy with  $\epsilon = \frac{1}{K}$  fulfills the GLIE property, where  $K$  is the episode counter
- thus,  $\epsilon$ -greedy policy will converge to optimal policy, with arbitrary starting states.

# On Policy vs. Off Policy

Exploration



Exploitation

On-policy methods  
need to find a compromise  
e.g. exploring starts,  $\epsilon$ -greedy, GLIE

"Learning on the job"

Off-policy methods

separate the problem

"Looking over the shoulder"

Behavior policy  $\pi_b$

generates behavior, data

Target policy  $\pi_T$

learns from  $\pi_b$ 's data

# Importance Sampling

Monte Carlo:  $\mathbb{E}[f(\underline{x})] \approx \hat{E} = \frac{1}{N} \sum_{i=1}^N f(\underline{x}_i)$ , where  $\underline{x}_i \stackrel{i.i.d.}{\sim} p_x$

Sometimes hard to sample  $p_x$ , but easy to sample a "proxy"  $q$ :

$$\begin{aligned}\mathbb{E}_p[f(\underline{x})] &= \int p_x(\underline{x}) f(\underline{x}) d\underline{x} \\ &= \int q(\underline{x}) \frac{p_x(\underline{x})}{q(\underline{x})} f(\underline{x}) d\underline{x} = \mathbb{E}_q \left[ \frac{p_x(\underline{x})}{q(\underline{x})} f(\underline{x}) \right]\end{aligned}$$

## Importance Sampling Estimator

$$\mathbb{E}[f(\underline{x})] \approx \hat{E}_{IS} = \frac{1}{N} \sum_{i=1}^N \frac{p_x(\underline{x}_i)}{q(\underline{x}_i)} f(\underline{x}_i), \text{ where } \underline{x}_i \stackrel{i.i.d.}{\sim} q$$

importance weights

# Off-policy Monte Carlo

- Generate an episode with behavior policy  $\pi_b$ :

$$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

- Probability of episode under target policy  $\pi$ :

$$p_{\pi}(A_0, R_1, \dots | S_0) = \prod_{t=0}^T p(S_t, R_t | S_0, A_0) \underbrace{\prod_{t=1}^{T-1} p(A_t | S_t)}_{MDP} \underbrace{p(S_T, R_T | S_{T-1}, A_{T-1})}_{MP}$$

- Probability of episode under behavior policy 6:

$$p_b(A_0, R_1, \dots | S_0) = b(A_0 | S_0) p(S_1, R_1 | S_0, A_0) b(A_1 | S_1) \dots p(S_T, R_T | S_{T-1}, A_{T-1})$$

- Importance weight for episode:

Importance weight for episode:

$$\frac{p_{\pi}(A_0, R_1, \dots | S_0)}{p_b(A_0, R_1, \dots | S_0)} = \frac{\prod_{t=0}^T p(A_t | S_t) \cancel{p(S_t, R_t | S_0, A_0)}}{\prod_{t=0}^T b(A_t | S_t) \cancel{p(S_t, R_t | S_0, A_0)}} \prod_{t=1}^T \frac{\prod_{s=0}^{t-1} p(S_s | S_0) \cancel{p(S_t, R_t | S_{t-1}, A_{t-1})}}{\prod_{s=0}^{t-1} b(S_s | S_0) \cancel{p(S_t, R_t | S_{t-1}, A_{t-1})}}$$

an known dynamics cancel!

$$= \prod_{t=0}^T \frac{\pi(A_t | S_t)}{g(A_t | S_t)}$$

## Off-policy Monte Carlo cont'd

- initialize  $q(s, a)$  arbitrarily  $\forall s \in S, a \in A$
- $N(s, a) \leftarrow 0 \quad \forall s \in S, a \in A$

- repeat

- generate episode using  $b$

$$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

- for  $t = T-1, \dots, 0$ :

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

$$W \leftarrow W \frac{\pi(A_t | S_t)}{b(A_t | S_t)} \quad // \text{recursive}$$

$$\prod_{t=t'}^T \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$$

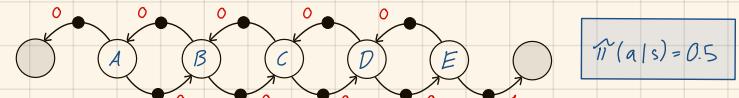
$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$q(S_t, A_t) \leftarrow q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (W q - q(S_t, A_t))$$

Correct, but silly algorithm... Why? (How is  $\pi$  defined?)

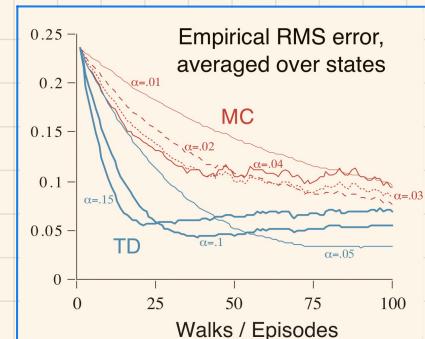
# Model-free Control using TD

# Sarsa: On-Policy TD Control



TD update for evaluation:

$$v(s) \leftarrow v(s) + \alpha (R_{t+1} + \gamma v(s_{t+1}) - v(s))$$

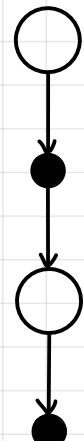
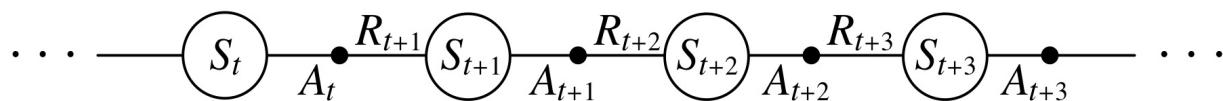


For control, we need TD update for  $q$ :

$$q(s,a) \leftarrow q(s,a) + \alpha (R_{t+1} + \gamma q(s_{t+1}, a_{t+1}) - q(s,a))$$

This rule is applied to episodes following current  $\hat{\pi}$  (on-policy),  
and to each consecutive quintuple  $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$

SARSA



# Sarsa: On-Policy TD Control

Let  $\epsilon > 0$ ,  $\alpha \in (0, 1]$  // hyperparameters for  $\epsilon$ -greedy, stepsize

Initialize  $q(s, a)$  arbitrarily, except  $q(s, a) = 0$  for terminal  $s$

- repeat: // for continuing tasks  $\Rightarrow$  one repetition

- init  $s$

- generate  $A$  from  $\epsilon$ -greedy( $q(s, \cdot)$ )

- repeat until  $s$  is terminal: // forever for continuing task

- take action  $A$

- observe  $s'$ ,  $R$

- generate  $A'$  from  $\epsilon$ -greedy( $q(s', \cdot)$ )

-  $q(s, A) \leftarrow q(s, A) + \alpha (R + \gamma q(s', A') - q(s, A))$

-  $s, A \leftarrow s', A'$

Sarsa converges to  $q^*$  and  $\pi^*$  if

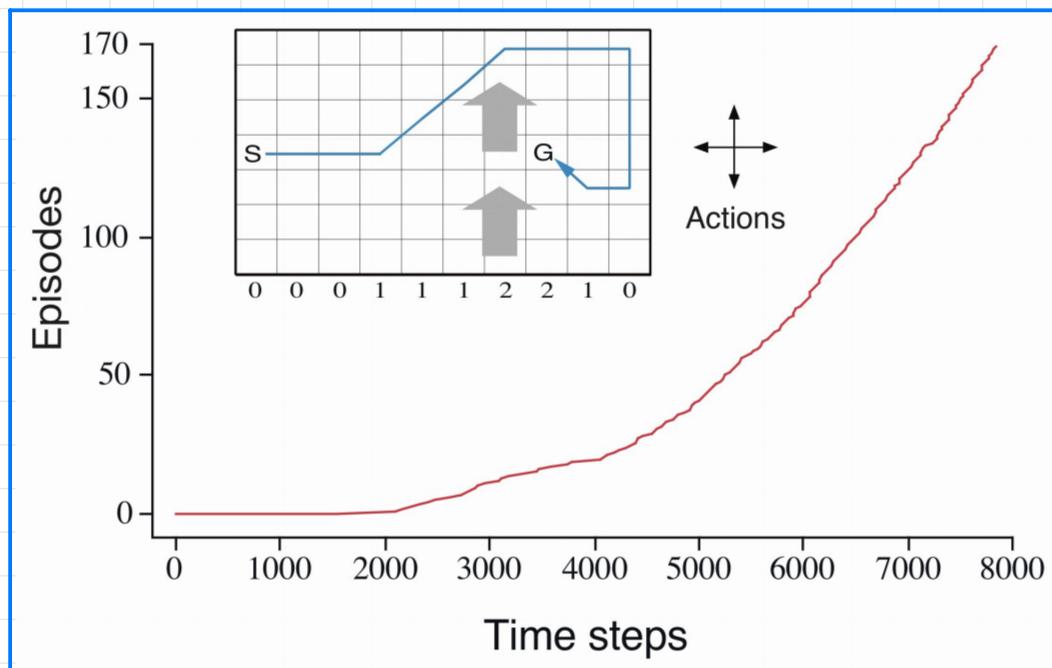
- GLIE holds

- $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$  (e.g.  $\alpha = \frac{1}{t}$ )

[Singh et al.,  
Machine Learning, 2000]

# Windy Grid World with Sarsa

- dedicated start state  $S$  and goal  $G$
- wind in the center, displacing agent by 1 or 2 cells
- Sarsa with  $\alpha = 0.5$  and  $\epsilon = 0.1$
- over time, agent learns to reach goal quicker



read: more completed episodes in fewer time steps

- Monte Carlo would learn very slowly, since it waits until end of episode before update — this takes very long in this grid world!

## Importance-sampled Sarsa

$$\hat{E}_{IS} = \frac{1}{N} \sum_{i=1}^N \frac{p_{\pi}(x_i)}{q(x_i)} f(x_i), \text{ where } x_i \stackrel{i.i.d.}{\sim} q$$

Sarsa:  $\hat{\pi} = \epsilon\text{-greedy}(q)$

$A_{t+1} \sim \hat{\pi}(a|s)$

$$q_t(s, a) \leftarrow q_t(s, a) + \alpha (R_{t+1} + \gamma q_t(S_{t+1}, A_{t+1}) - q_t(s, a))$$

- On-policy,  $\hat{\pi}$  is learned while following it
- Using importance sampling, we can make Sarsa off-policy

## Importance-weighted Sarsa:

$b = \epsilon\text{-greedy}(q)$  // behavior

$\hat{\pi} = \text{greedy}(q)$  // target

$A_{t+1} \sim b(a|s)$

$$q(s, a) \leftarrow q(s, a) + \alpha \left( R_{t+1} + \gamma \frac{\hat{\pi}(A_{t+1}|S_{t+1})}{b(A_{t+1}|S_{t+1})} q(S_{t+1}, A_{t+1}) - q(s, a) \right)$$

- no GLE required,  $\epsilon$  can be constant
- has rather large variance  $\rightarrow$  Q-learning generally preferred

# Q-learning: Off-Policy TD Control

[Watkins 1989]

Sarsa:  $q_t(s, a) \leftarrow q_t(s, a) + \alpha \underbrace{(R_{t+1} + \gamma q_t(S_{t+1}, A_{t+1}) - q_t(s, a))}_{\mathbb{E}[\delta_t] = 0 \triangleq \text{Bellman expectation equation}}$

Q-learning  $q(s, a) \leftarrow q(s, a) + \alpha \underbrace{(R_{t+1} + \gamma \max_{a'} q(S_{t+1}, a') - q(s, a))}_{\mathbb{E}[\delta_t] = 0 \triangleq \text{Bellman optimality equation}}$

Let  $\epsilon > 0$ ,  $\alpha \in (0, 1]$  // hyperparameters for  $\epsilon$ -greedy, stepsize

Initialize  $q(s, a)$  arbitrarily, except  $q(s, a) = 0$  for terminal  $s$

- repeat:

- init  $S$

- repeat until  $S$  is terminal:

- take action  $A$  from behavioral policy, e.g.  $\epsilon$ -greedy ( $q(S, \cdot)$ )

- observe  $S'$ ,  $R$

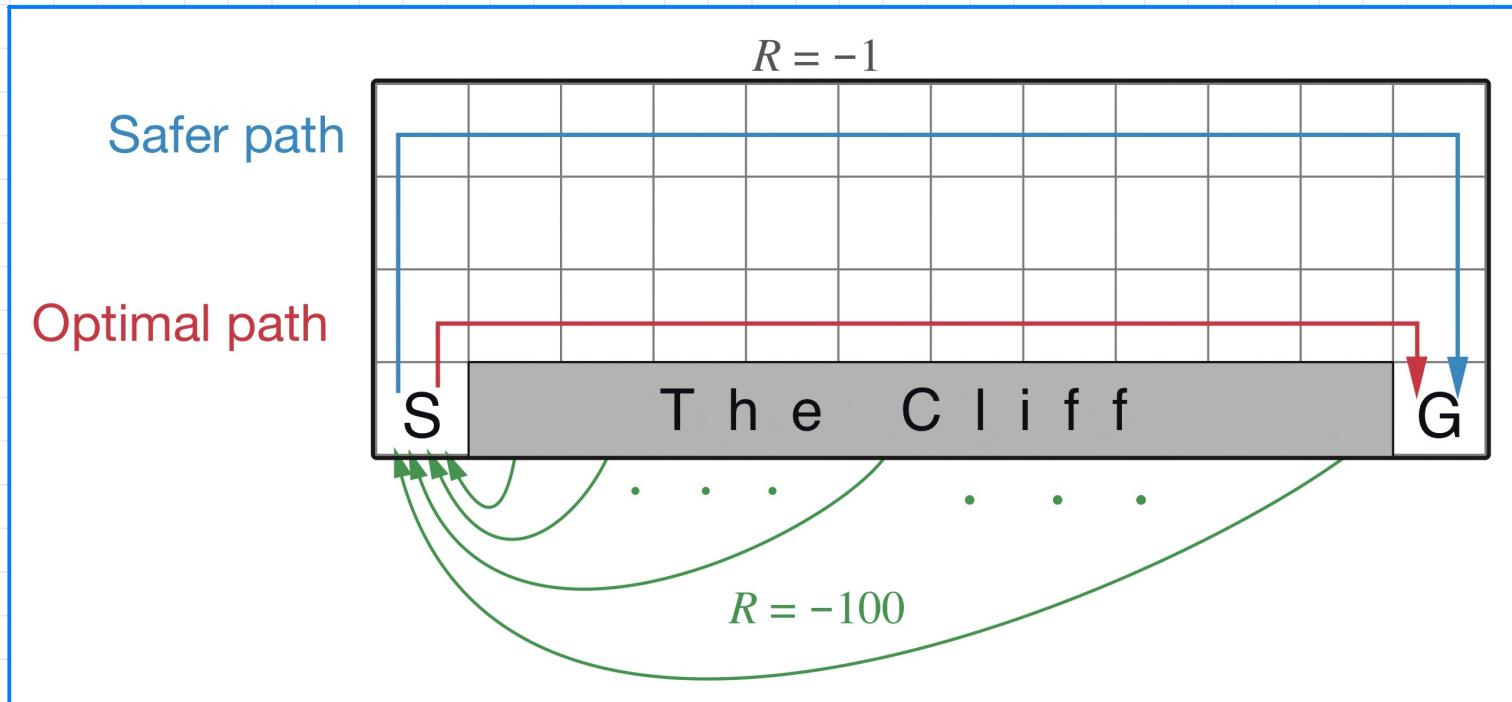
-  $q(S, A) \leftarrow q(S, A) + \alpha \left( R + \gamma \max_{a'} q(S', a') - q(S, A) \right)$

-  $S \leftarrow S'$

Off-policy: learn (hard) greedy policy while following  $\epsilon$ -greedy policy.  
 $\epsilon$  can be constant (no GLIE). Convergence to  $q^*$  for  $\sum_k \alpha_k = \infty$ ,  $\sum_k \alpha_k^2 < \infty$ .

# The Cliff

- grid world with starting state and goal
- the grey area (the cliff) yields large penalty ( $R=-100$ ) and transition to the start
- when running Sarsa and Q-learning with constant  $\epsilon=0.1$ , one will learn the optimal path and the other the safer path.
- which algorithm will learn which path?



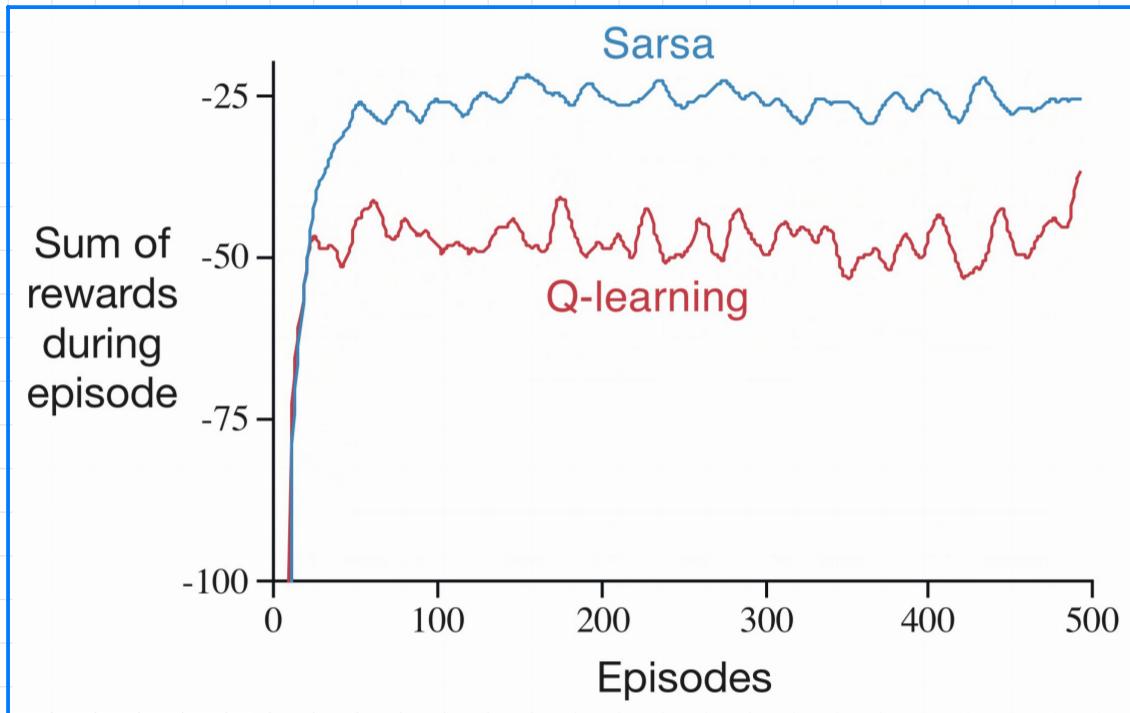
## The Cliff cont'd

Q-learning: optimal path

Sarsa: safer path

Sarsa is on-policy: return of behavior policy is maximized

Q-learning is off-policy: return of behavior policy not a goal per se



If  $\epsilon = 1/\kappa$  (GLIE)

- Sarsa converges to  $\pi^*$
- Q-learning achieves optimal online returns

## Expected Sarsa

greedy policy:  $\pi'(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q(s, a') \\ 0 & \text{otherwise} \end{cases}$

Q-Learning:  $q(s, a) \leftarrow q(s, a) + \alpha (R_{t+1} + \gamma \max_{a'} q(S_{t+1}, a') - q(s, a))$

- off-policy
- behavior  $\epsilon$ -greedy
- target greedy

$$E_{\pi'} [q(S_{t+1}, A_{t+1})] = \sum_{a'} \pi'(a' | S_{t+1}) q(S_{t+1}, a')$$

## Sarsa:

$$q(s, a) \leftarrow q(s, a) + \alpha (R_{t+1} + \gamma q(S_{t+1}, A_{t+1}) - q(s, a))$$

- on-policy
- $\epsilon$ -greedy (GLIE)

$$\hat{E} \approx E_{\pi} [q(S_{t+1}, A_{t+1})], \quad A_{t+1} \sim \pi(a|s)$$

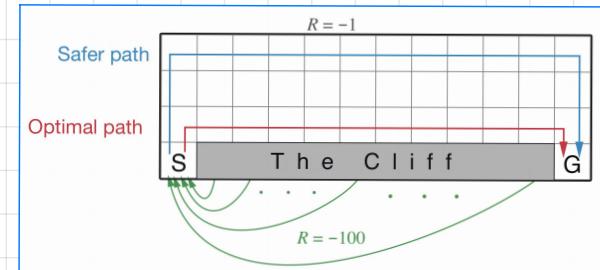
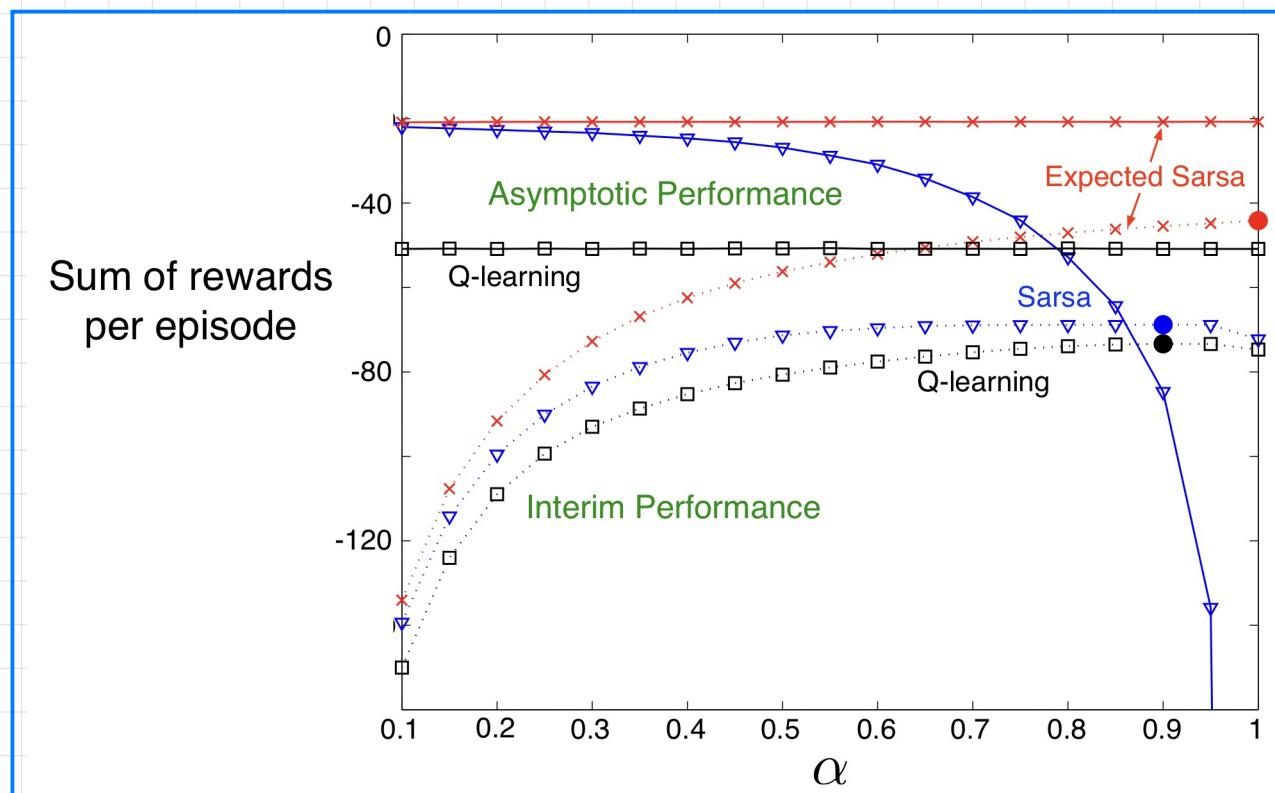
(single sample Monte Carlo)

## Expected Sarsa

$$q(s, a) \leftarrow q(s, a) + \alpha (R_{t+1} + \gamma \sum_{a'} \pi'(a' | S_{t+1}) q(S_{t+1}, a') - q(s, a))$$

- on-policy
- $\epsilon$ -greedy (GLIE)
- lower variance than Sarsa
- turns into Q-learning when using  $\pi'$  instead

# The Cliff cont'd



- all algorithms used  $\epsilon = 0.1$
- interim performance: 100 episodes
- asymptotic performance: 100,000 episodes
- expected Sarsa learns faster than Sarsa  
also less sensitive to choice of  $\alpha$

# Summary TD Control

on-policy

Sarsa

sample A



off-policy

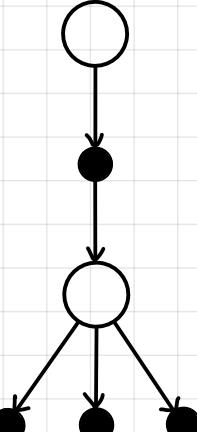
Importance-weighted Sarsa



$$\frac{\hat{\pi}(A_{t+1} | S_{t+1})}{b(A_{t+1} | S_{t+1})}$$

Expected Sarsa

$E_A$



Q-learning (Sarsa Max)

