

ST455: Reinforcement Learning

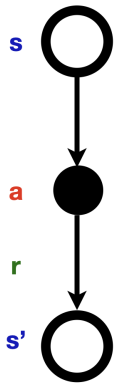
Lecture 4: Temporal Difference (TD) Learning

Chengchun Shi

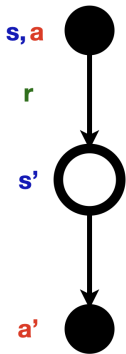
Lecture Outline

1. TD Prediction
2. SARSA
3. Q-Learning
4. TD(λ) and SARSA(λ)

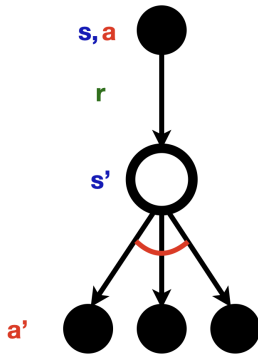
Lecture Outline (Cont'd)



TD

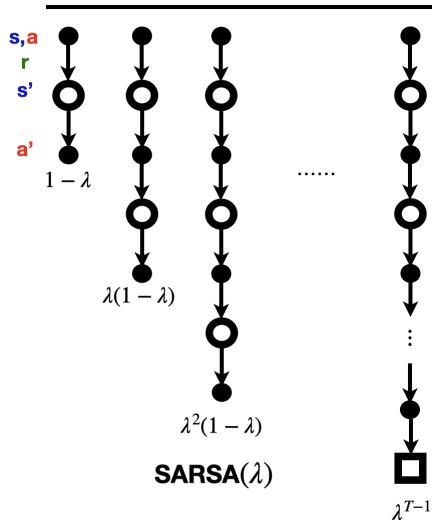
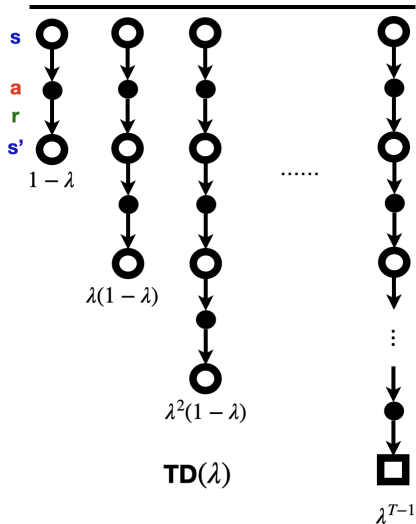


SARSA



Q-Learning

Lecture Outline (Cont'd)



1. TD Prediction

2. SARSA

3. Q-Learning

4. $TD(\lambda)$ and $SARSA(\lambda)$

TD Learning v.s. MC Methods

- **Temporal Difference (TD) Learning**: a learning method that combines ideas from Monte Carlo (MC) methods and dynamic programming
- **Similarities** with MC methods:
 - TD is also **model free**, i.e. learns directly from the experience **without** a model (e.g., state transitions, reward functions)
 - Difference from dynamic programming
- **Difference** with MC methods:
 - MC methods: updates the value **episode-by-episode** (applies to **episodic** tasks)
 - TD methods: updates the value **step-by-step** (applies to **continuous** tasks as well)

MC Prediction

- **Objective:** learns V^π from experience under π
- MC Policy Evaluation: $V(s) \leftarrow \text{average}[\text{Returns}(s)]$
- Incremental update for every-visit MC prediction:

$$V(S_t) \leftarrow V(S_t) + \alpha_t [G_t - V(S_t)]$$

where α_t is $\frac{1}{\#[\text{Returns}(S_t)]}$ at time t

- We may regard G_t as a **target**
- The update can be performed after return G_t is observed
- i.e. after the episode is completed

TD Prediction

- Unlike MC methods, TD methods wait only until **next time step**
- The simplest TD method (so called TD(0)) considers the update

$$V(S_t) \leftarrow V(S_t) + \alpha_t [R_t + \gamma V(S_{t+1}) - V(S_t)]$$

- This update rule has $R_t + \gamma V(S_{t+1})$ as the **target**
- Considered as a **bootstrap** method: update in part based on an existing estimate
- Different from “bootstrap” in statistics: a resampling method (e.g., sample with replacement) for **uncertainty quantification** of a given estimate

MC vs TD update

- Notice that under the MDP assumption

$$\mathbf{V}^\pi(\mathbf{s}) = \mathbb{E}^\pi(\mathbf{G}_t | \mathbf{S}_t = \mathbf{s}) \quad (1)$$

$$\begin{aligned} &= \mathbb{E}^\pi\left(\sum_{k=0}^{\infty} \gamma^k \mathbf{R}_{t+k} | \mathbf{S}_t = \mathbf{s}\right) \\ &= \mathbb{E}^\pi[\mathbf{R}_t + \gamma \mathbf{V}^\pi(\mathbf{S}_{t+1}) | \mathbf{S}_t = \mathbf{s}] \end{aligned} \quad (2)$$

- MC methods use as the target the random variable in (1)
- TD methods use as the target the random variable in (2)
 - **Immediate reward** and estimate of the **future value**

Bootstrapping and Sampling

- **Bootstrapping:** update involves an estimate
 - MC does **not** bootstrap
 - DP bootstraps
 - TD bootstraps
- **Sampling:** update samples an expectation
 - MC samples
 - DP does **not** sample
 - TD samples

TD(0): Pseudocode

- **Input:** π policy to be evaluated, step size α
- **Initialization:** V arbitrary
- **Repeat** for each episode:

Initialize state s

Repeat for each step of the episode:

$a \leftarrow$ action given by π for s

 Take action a , observe reward r and next state s'

$V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$

$s \leftarrow s'$

until s is a terminal state

Pros & Cons of MC vs TD

- MC must wait until the **end** of episode
- MC learns from **complete** sequences
- MC only works for **episodic** (terminating) environments
- TD can learn online after **each step**
- TD can learn from **incomplete** sequences
- TD works in **continuing** environments

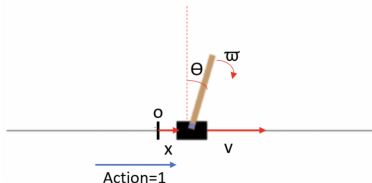
Pros & Cons of MC vs TD (Cont'd)

- Bias/Variance Trade-Off
- Return G_t is **unbiased** estimate of $V^\pi(S_t)$
- Oracle target $R_t + \gamma V^\pi(S_{t+1})$ is **unbiased** estimate of $V^\pi(S_t)$
- TD target $R_t + \gamma V(S_{t+1})$ is **biased** estimate of $V^\pi(S_t)$
- TD target has much lower **variance** than the return
 - Return depends on **many** random actions, transitions, rewards
 - TD target depends on **one** random action, transition, reward
- MC has **high variance, zero bias, insensitive to initialization**
- TD has **low variance, some bias, sensitive to initialization**

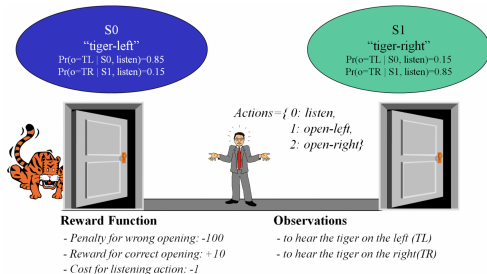
Pros & Cons of MC vs TD (Cont'd)

- TD exploits **Markov** & **stationary** properties
- Relies on the **Bellman equation**
- More **efficient** in MDP environments

frame: 53, Obs: (0.018, 0.669, 0.286, 0.618)
Action: 1.0, Cumulative Reward: 47.0, Done: 1



- MC **does** not exploit these properties
- More **flexible** in non-MDP environments (e.g., POMDP)



Rate of Convergence

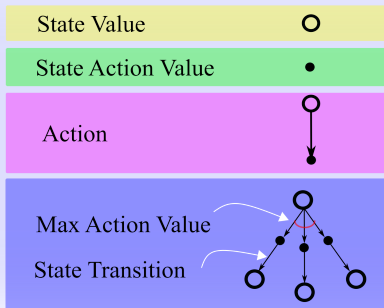
- For i.i.d. random variables $\mathbf{X}_1, \dots, \mathbf{X}_n$ with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\sigma}^2$,

$$\sqrt{n}(\bar{\mathbf{X}} - \boldsymbol{\mu}) \rightarrow N(\mathbf{0}, \boldsymbol{\sigma}^2),$$

according to CLT.

- $\bar{\mathbf{X}}$ converges to $\boldsymbol{\mu}$ at a rate of $n^{-1/2}$.
- For n episodes with T time points per episode, first-visit MC converges at a rate of $n^{-1/2}$.
- For n episodes with T time points per episode, TD converges at a rate of $(nT)^{-1/2}$, with proper choice of step sizes [see e.g., Tadić, 2002].
- First-visit MC requires $n \rightarrow \infty$ to be consistent
- TD requires either n or $T \rightarrow \infty$ to be consistent

Backup Diagram



Taken from <https://towardsdatascience.com/all-about-backup-diagram-fefb25aaf804>

Backup Diagram (Cont'd)



TD



More states
and actions

Terminal state

MC

1. TD Prediction

2. SARSA

3. Q-Learning

4. $TD(\lambda)$ and $SARSA(\lambda)$

SARSA: an On-Policy TD Control

- **SARSA**: a TD method for **policy optimisation**
 - Follows the pattern of policy iteration
 - Uses TD prediction method for **policy evaluation**
 - Uses ϵ -greedy exploration for **policy improvement**
- Similar to MC control, estimate state-action value $Q^\pi(s, a)$ (instead of the state value $V^\pi(s)$) for the control problem
- Different from MC control, update the state-value **every time step**

Bellman Equations

- Bellman equation for the (state) value function:

$$V^\pi(s) = \mathbb{E}[R_t + \gamma V^\pi(S_{t+1}) | S_t = s].$$

- Bellman equation for the state-action value function:

$$Q^\pi(s, a) = \mathbb{E} \left[R_t + \gamma \sum_{a'} \pi(a' | S_{t+1}) Q^\pi(S_{t+1}, a') | A_t = a, S_t = s \right],$$

or equivalently,

$$Q^\pi(s, a) = \mathbb{E}^\pi[R_t + \gamma Q^\pi(S_{t+1}, A_{t+1}) | A_t = a, S_t = s].$$

SARSA: Policy Evaluation

- Incremental estimation of the state-action value function:

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

for non-terminal state S_{t+1}

- If S_{t+1} is a terminal state, $Q(S_{t+1}, A_{t+1}) = 0$
- This update uses every element of the quintuple of variables:

$$(S_t, A_t, R_t, S_{t+1}, A_{t+1})$$

$S \quad A \quad R \quad S \quad A$

SARSA: Pseudocode

- **Initialization:** Q arbitrary
- **Repeat** for each episode:
 - Initialize** state s
 - Choose** action a from s using policy derived from Q (ϵ -greedy)
 - Repeat** for each step of the episode:
 - Take action a , observe reward r and next state s'
 - $a' \leftarrow$ action from s' using policy derived from Q (ϵ -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 a terminal state

Convergence of SARSA

Theorem

SARSA converges to the optimal Q-function, $Q(\mathbf{s}, \mathbf{a}) \rightarrow Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})$ for any \mathbf{s} and \mathbf{a} , if

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

$$\sum_{t=0}^{\infty} \mathbb{I}(\mathbf{S}_t = \mathbf{s}, \mathbf{A}_t = \mathbf{a}) = \infty.$$

- The policy converges to a greedy policy,

$$\lim_{t \rightarrow \infty} \pi_t(\mathbf{a} | \mathbf{s}) = \mathbb{I}(\mathbf{a} = \arg \max_{\mathbf{a}'} Q_t(\mathbf{s}, \mathbf{a}'))$$

- Robbins-Monro sequence of step-sizes [Robbins and Monroe, 1951],

$$\sum_{t=0}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

Convergence of SARSA (Cont'd)

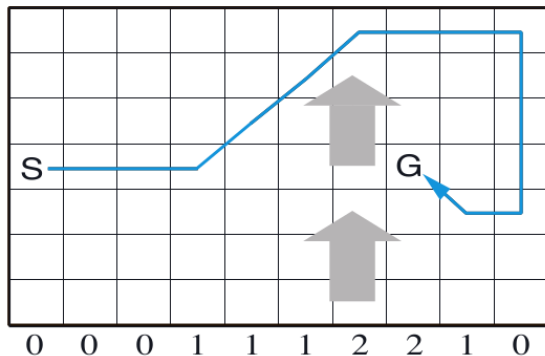
- Condition 1: All state-action pairs are **explored** infinitely many times
 $\Rightarrow \epsilon$ to be strictly positive
- Condition 2: The policy converges to a **greedy** policy
 $\Rightarrow \epsilon_t$ decays to zero as t grows to infinity
- Condition 3: Robbins-Monro sequence of step-sizes

$$\sum_{t=0}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

$\Rightarrow \alpha_t$ proportional to t^{-c} for $1/2 < c \leq 1$.

$\sum_t t^{-c} = \infty$ when $c \leq 1$ and $\sum_t t^{-c} < \infty$ when $c > 1$

Windy Gridworld Example



- An **episodic** task
- Rewards of -1 until **goal** is reached
- Strength of wind indicated by numbers

Windy Gridworld Example (Cont'd)

Optimal policy is:

['R'	'D'	'R'	'R'	'R'	'R'	'R'	'R'	'R'	'D']
['R'	'R'	'U'	'L'	'R'	'R'	'R'	'U'	'R'	'D']
['L'	'U'	'R'	'R'	'R'	'R'	'R'	'R'	'R'	'D']
['R'	'R'	'R'	'R'	'R'	'R'	'U'	'G'	'R'	'D']
['R'	'D'	'D'	'R'	'R'	'R'	'U'	'D'	'L'	'L']
['D'	'D'	'D'	'R'	'R'	'U'	'U'	'D'	'R'	'U']
['R'	'R'	'U'	'R'	'U'	'U'	'U'	'U'	'R'	'L']
['0'	'0'	'0'	'1'	'1'	'1'	'2'	'2'	'1'	'0']

1. TD Prediction

2. SARSA

3. Q-Learning

4. $TD(\lambda)$ and $SARSA(\lambda)$

Q-Learning

- One of the most **popular** class of RL algorithms
- Variants include double Q-learning, fitted Q-iteration, deep Q-network (DQN), quantile DQN (more in later lectures)
- **Main idea**: learn the optimal Q-function $Q^{\pi^{\text{opt}}}$ based on the Bellman **optimality** equation and derive the optimal policy (see Appendix for the proof)

$$\pi^{\text{opt}}(\textcolor{blue}{s}) = \arg \max_{\textcolor{red}{a}} Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})$$

- Focus on **tabular Q-learning** in this lecture (finite MDP, discrete state and action)

Bellman Optimality Equation

- Bellman optimality equation for the optimal value function:

$$V^{\pi^{\text{opt}}}(s) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(S_{t+1}) | A_t = a, S_t = s].$$

- Bellman optimal equation for the optimal Q-function (see Appendix for the proof):

$$Q^{\pi^{\text{opt}}}(s, a) = \mathbb{E} \left[R_t + \gamma \max_{a'} Q^{\pi^{\text{opt}}}(S_{t+1}, a') | A_t = a, S_t = s \right].$$

Q-Learning: an Off-Policy TD Control

- One-step SARSA update:

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha[\mathbf{R}_t + \gamma Q(\mathbf{S}_{t+1}, \mathbf{A}_{t+1}) - Q(\mathbf{S}_t, \mathbf{A}_t)]$$

- One-step Q-learning update:

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha \left[\mathbf{R}_t + \gamma \max_a Q(\mathbf{S}_{t+1}, a) - Q(\mathbf{S}_t, \mathbf{A}_t) \right]$$

- In Q-learning, the action in the target is independent of the **behavior policy**
- The behavior policy has an effect on which state-actions are visited

Q-Learning: Pseudocode

- **Initialization:** Q arbitrary
- **Repeat** for each episode:
 - Initialize** state s
 - Repeat** for each step of the episode:
 - $a \leftarrow$ action from s using policy derived from Q (e.g., ϵ -greedy)
 - Take action a , observe reward r and next state 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 a terminal state

On-Policy v.s Off-Policy

- Q-learning is **off-policy**:
 - Updates Q-values using Q-value of next state s' and **greedy** action a'
 - Assumes greedy policy were followed despite that it's not following greedy policy
- SARSA is **on-policy**:
 - Updates Q-values using Q-value of next state s' and **current policy's** action a'
 - Assumes the current policy continues to be followed

Convergence of Q-Learning

Theorem (Melo [2001])

Q-learning converges to the optimal Q-function if

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

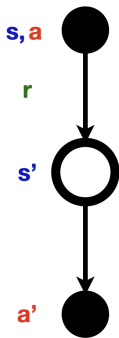
$$\sum_{t=0}^{\infty} \mathbb{I}(\mathbf{S}_t = \mathbf{s}, \mathbf{A}_t = \mathbf{a}) = \infty.$$

- *Robbins-Monro sequence of step-sizes [Robbins and Monro, 1951],*

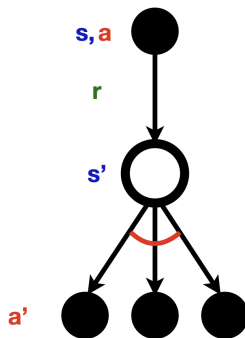
$$\sum_{t=0}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

- Only requires ε to be strictly positive
- **No** need to require ε to decay to zero
- Q-learning converges even if the behavior policy is far from the optimal

Backup Diagram



SARSA



Q-Learning

Cliff Walking Example

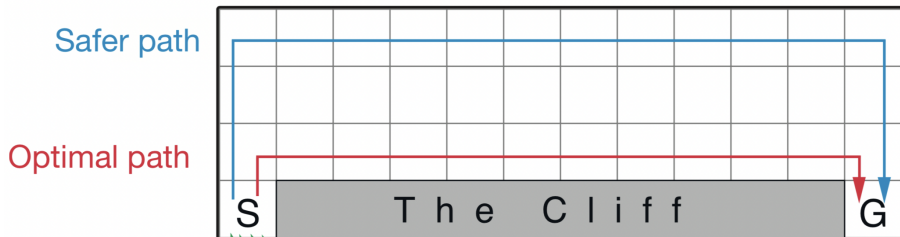


Figure: Illustrations of Cliff Walking

- Undiscounted, **episodic** task
- Actions: up, down, right and left
- Reward of -100 if stepping into **cliff**
- Reward of -1 on other transitions

Cliff Walking Example (Cont'd)

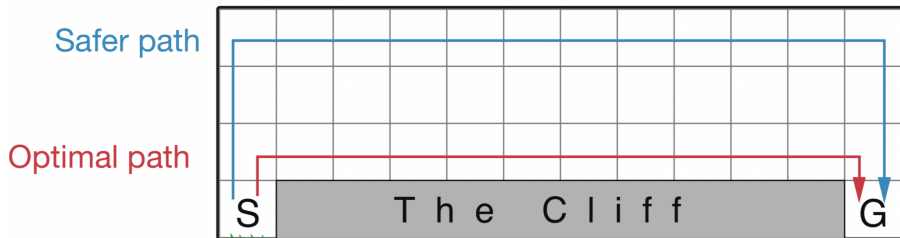


Figure: Illustrations of Cliff Walking

- Q-learning identifies the **optimal** path
- SARSA identifies a **safer** path (the optimal path is not optimal here due to that the ϵ -greedy policy, which might force the agent to fall into the cliff when walking along the optimal path, yielding a low value)

Maximization Bias

- One-step Q-learning update:

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha \left[R_t + \gamma \max_{\mathbf{a}} Q(\mathbf{S}_{t+1}, \mathbf{a}) - Q(\mathbf{S}_t, \mathbf{A}_t) \right],$$

- Maximum over $Q(\mathbf{S}_{t+1}, \mathbf{a})$ can lead to **significant positive** bias
- **Example:**
 - Oracle optimal Q-function $Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a}) = \mathbf{0}$ for any (\mathbf{s}, \mathbf{a})
 - $\max_{\mathbf{a}} Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a}) = \mathbf{0}$ for any \mathbf{s}
 - Estimated Q-function $Q(\mathbf{s}, \mathbf{a})$: uncertain, some above and some below zero
 - $\max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ likely to be **positive**

Maximization Bias (Cont'd)

- Maximization over Q involves two steps: **greedy-action selection** and **state-action value evaluation**

$$\max_a Q(s, a) = Q(s, \arg \max_a Q(s, a))$$

- Solution: use two different Q -functions for two steps

$$Q_1(s, \arg \max_a Q_2(s, a))$$

- **Example:** $Q^{\pi^{\text{opt}}} = 0$. Due to difference between Q_1 and Q_2 , the above expression is no longer always positive

Double Q-Learning

- **Initialize** two Q-functions Q_1 and Q_2
- **Divide** time steps into two by flipping a coin on each step
- If the coin comes up with head

$$Q_1(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q_1(\mathbf{S}_t, \mathbf{A}_t) + \alpha \left[R_t + \gamma Q_2(\mathbf{S}_{t+1}, \arg \max_{\mathbf{a}} Q_1(\mathbf{S}_{t+1}, \mathbf{a})) - Q_1(\mathbf{S}_t, \mathbf{A}_t) \right]$$

- Otherwise

$$Q_2(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q_2(\mathbf{S}_t, \mathbf{A}_t) + \alpha \left[R_t + \gamma Q_1(\mathbf{S}_{t+1}, \arg \max_{\mathbf{a}} Q_2(\mathbf{S}_{t+1}, \mathbf{a})) - Q_2(\mathbf{S}_t, \mathbf{A}_t) \right]$$

Double Q-Learning: Pseudocode

- **Initialization:** Q_1 and Q_2 arbitrary

- **Repeat** for each episode:

Initialize state s

Repeat for each step of the episode:

$a \leftarrow$ action from s using policy derived from $Q_1 + Q_2$ (e.g., ϵ -greedy)

 Take action a , observe reward r and next state s'

 With probability **0.5**:

$a' \leftarrow \arg \max_a Q_1(s_{t+1}, a)$

$Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha [R_t + \gamma Q_2(s_{t+1}, a') - Q_1(s_t, a_t)]$

 else:

$a' \leftarrow \arg \max_a Q_2(s_{t+1}, a)$

$Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha [R_t + \gamma Q_1(s_{t+1}, a') - Q_2(s_t, a_t)]$

$s \leftarrow s'$

until s is a terminal state

1. TD Prediction

2. SARSA

3. Q-Learning

4. $TD(\lambda)$ and $SARSA(\lambda)$

n -Step Return

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

$$n = 1 \quad (\text{TD}) \quad G_t^{(1)} = R_t + \gamma V^\pi(S_{t+1})$$

$$n = 2 \quad G_t^{(2)} = R_t + \gamma R_{t+1} + \gamma^2 V^\pi(S_{t+2})$$

$$\vdots \quad \vdots$$

$$n = \infty \quad (\text{MC}) \quad G_t^{(\infty)} = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

- Define n -step return

$$G_t^{(n)} = R_t + \gamma R_{t+1} + \dots + \gamma^n V^\pi(S_{t+n})$$

- n -step temporal difference learning

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t^{(n)} - V(S_t)]$$

Averaging n -Step Return

- We can average n -step returns over different n
- e.g. average the 2-step and 4-step returns

$$\frac{1}{2}G_t^{(2)} + \frac{1}{2}G_t^{(4)}$$

- Combines information from two different time-steps
- Can we combine information from all time-steps?

λ -Return

- The λ -return G_t^λ combines all n -step return $G_t^{(n)}$
- Using weight $(1 - \lambda)\lambda^{n-1}$

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

- Notice that

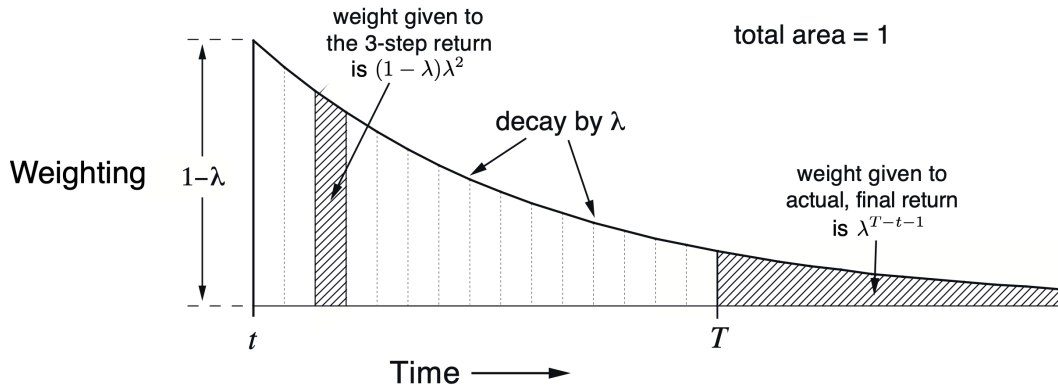
$$\sum_{n=1}^{+\infty} (1 - \lambda) \lambda^{n-1} = 1$$

- TD(λ)

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t^\lambda - V(S_t)]$$

- Like MC, can only be computed from complete episodes

Weighting Function



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

Special Cases

- TD(λ)

$$V(\mathbf{S}_t) \leftarrow V(\mathbf{S}_t) + \alpha[\mathbf{G}_t^\lambda - V(\mathbf{S}_t)]$$

- When $\lambda = 0$, reduces to TD method

$$\begin{aligned} V(\mathbf{S}_t) &\leftarrow V(\mathbf{S}_t) + \alpha[\mathbf{G}_t^{(1)} - V(\mathbf{S}_t)] \\ &= V(\mathbf{S}_t) + \alpha[\mathbf{R}_t + \gamma V(\mathbf{S}_{t+1}) - V(\mathbf{S}_t)] \end{aligned}$$

- When $\lambda = 1$, reduces to MC method

$$\begin{aligned} V(\mathbf{S}_t) &\leftarrow V(\mathbf{S}_t) + \alpha[\mathbf{G}_t^{(1)} - V(\mathbf{S}_t)] \\ &= V(\mathbf{S}_t) + \alpha[\mathbf{R}_t + \gamma \mathbf{R}_{t+1} + \dots + \gamma^T \mathbf{R}_{t+T} - V(\mathbf{S}_t)] \end{aligned}$$

n -Step SARSA

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

$$n = 1 \quad (\text{SARSA}) \quad Q_t^{(1)} = R_t + \gamma Q^\pi(S_{t+1}, A_{t+1})$$

$$n = 2 \quad Q_t^{(2)} = R_t + \gamma R_{t+1} + \gamma^2 Q^\pi(S_{t+2}, A_{t+2})$$

$$\vdots \quad \quad \quad \vdots$$

$$n = \infty \quad (\text{MC}) \quad Q_t^{(\infty)} = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

- Define n -step return

$$Q_t^{(n)} = R_t + \gamma R_{t+1} + \dots + \gamma^n Q^\pi(S_{t+n}, A_{t+n})$$

- n -step Sarsa updates

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [Q_t^{(n)} - Q(S_t, A_t)]$$

SARSA(λ)

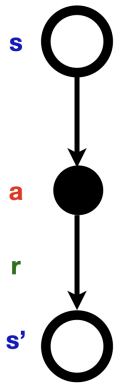
- The Q^λ -return combines all n -step return $Q_t^{(n)}$
- Using weight $(1 - \lambda)\lambda^{n-1}$

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

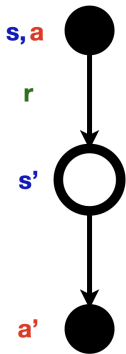
- SARSA(λ)

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha [Q_t^\lambda - Q(\mathbf{S}_t, \mathbf{A}_t)]$$

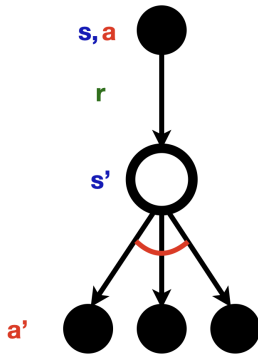
Summary



TD

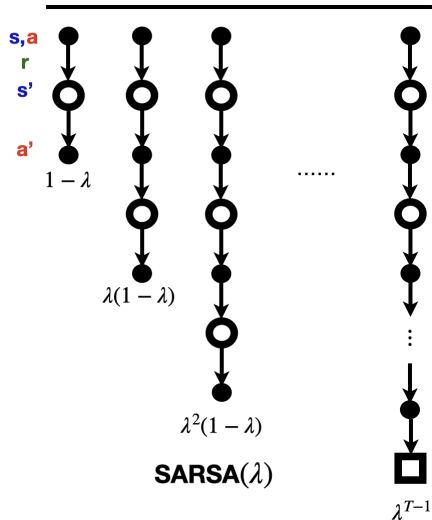
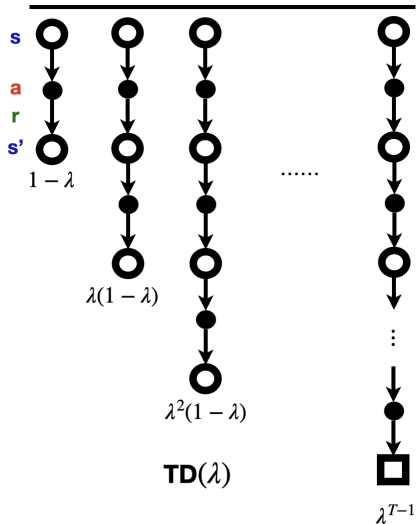


SARSA



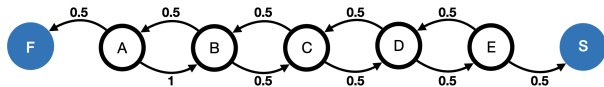
Q-Learning

Summary (Cont'd)

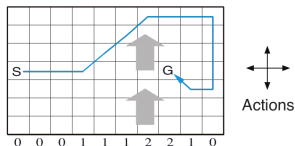


Seminar Exercises

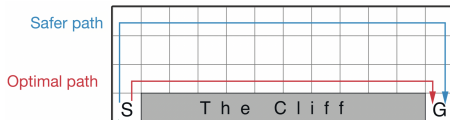
- Solution to HW3 (Deadline: Wed 12pm)
- TD: Random Walk



- Sarsa: Windy GridWorld



- Q-Learning: Cliff Walking Example



References I

- Francisco S Melo. Convergence of q-learning: A simple proof. *Institute Of Systems and Robotics, Tech. Rep*, pages 1–4, 2001.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- Vladislav B Tadić. On the almost sure rate of convergence of temporal-difference learning algorithms. *IFAC Proceedings Volumes*, 35(1):455–460, 2002.

Questions

Appendix: $\pi^{\text{opt}}(s) = \arg \max_a Q^{\pi^{\text{opt}}}(s, a)$?

- $Q^{\pi^{\text{opt}}}(s, a)$ is the **value** of the policy that
 - Assigns a at the initial decision time;
 - Follow π^{opt} afterwards
- $Q^{\pi^{\text{opt}}}(s, \pi^{\text{opt}}(s)) = V^{\pi^{\text{opt}}}(s)$ is the value under the optimal policy π^{opt}
- π^{opt} is stationary and is no worse than any **history-dependent** policies (Lecture 2)

$$Q^{\pi^{\text{opt}}}(s, a) \leq V^{\pi^{\text{opt}}}(s) = Q^{\pi^{\text{opt}}}(s, \pi^{\text{opt}}(s)), \quad \forall a.$$

- It follows that

$$\pi^{\text{opt}}(s) = \arg \max_a Q^{\pi^{\text{opt}}}(s, a)$$

Appendix: Proof of Bellman Optimality Equation

- Bellman optimal equation for the optimal Q-function:

$$Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a}) = \mathbb{E} \left[\textcolor{green}{R}_t + \textcolor{orange}{\gamma} \max_{\textcolor{red}{a}'} Q^{\pi^{\text{opt}}}(\textcolor{blue}{S}_{t+1}, \textcolor{red}{a}') | \textcolor{red}{A}_t = \textcolor{red}{a}, \textcolor{blue}{S}_t = \textcolor{blue}{s} \right].$$

- **Proof:** according to Bellman equation,

$$Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a}) = \mathbb{E} \left[\textcolor{green}{R}_t + \textcolor{orange}{\gamma} Q^{\pi^{\text{opt}}}(\textcolor{blue}{S}_{t+1}, \pi^{\text{opt}}(\textcolor{blue}{S}_{t+1})) | \textcolor{red}{A}_t = \textcolor{red}{a}, \textcolor{blue}{S}_t = \textcolor{blue}{s} \right]$$

- Since $\pi^{\text{opt}}(\textcolor{blue}{s}) = \arg \max_{\textcolor{red}{a}} Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})$, it follows that

$$\max_{\textcolor{red}{a}'} Q^{\pi^{\text{opt}}}(\textcolor{blue}{S}_{t+1}, \textcolor{red}{a}') = Q^{\pi^{\text{opt}}}(\textcolor{blue}{S}_{t+1}, \pi^{\text{opt}}(\textcolor{blue}{S}_{t+1}))$$