

# Model-Free Tabular Method Solutions

Alina Vereshchaka

CSE4/510 Reinforcement Learning  
Fall 2019

*avereshc@buffalo.edu*

September 17, 2019

# Overview

1 Recap

2 SARSA

# Table of Contents

1 Recap

2 SARSA

# Recap: Monte Carlo (MC) and Temporal Difference (TD) Learning

- ▶ **Goal:** learn  $v_\pi(s)$  from episodes of experience under policy  $\pi$

- ▶ Incremental **every-visit Monte-Carlo**:

- Update value  $V(S_t)$  toward actual return  $G_t$

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

- ▶ Simplest **Temporal-Difference** learning algorithm: TD(0)

- Update value  $V(S_t)$  toward estimated returns  $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

- ▶  $R_{t+1} + \gamma V(S_{t+1})$  is called the **TD target**
- ▶  $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is called the **TD error**.

# Recap: DP vs. MC vs TD Learning

► Remember:

MC: sample average return  
approximates expectation

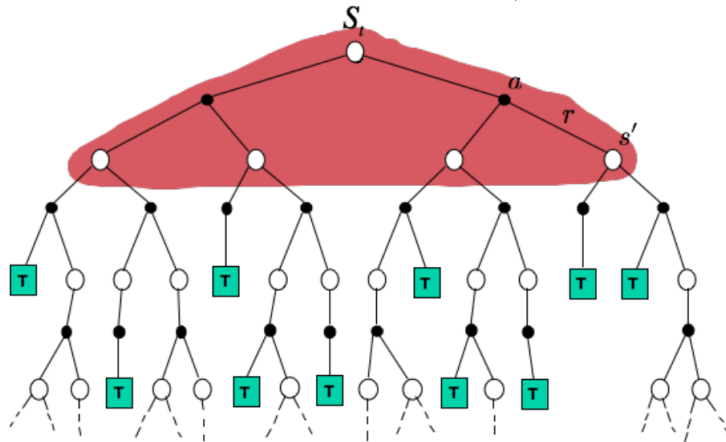
$$\begin{aligned}v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\&= \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right] \\&= \mathbb{E}_{\pi}\left[R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \mid S_t = s\right] \\&= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s].\end{aligned}$$

TD: combine both: Sample  
expected values and use a  
current estimate  $V(S_{t+1})$  of the  
true  $v_{\pi}(S_{t+1})$

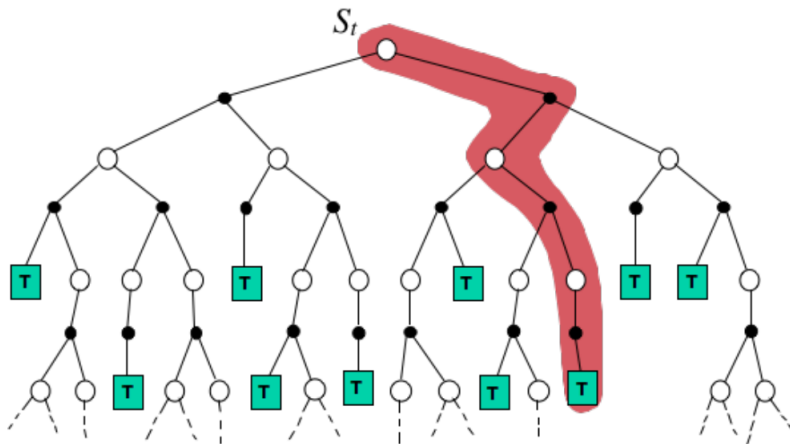
DP: the expected values are  
provided by a model. But we  
use a current estimate  $V(S_{t+1})$   
of the true  $v_{\pi}(S_{t+1})$

# Dynamic Programming

$$V(S_t) \leftarrow E_\pi[R_{t+1} + \gamma V(S_{t+1})] = \sum_a \pi(a|S_t) \sum_{s', r} p(s', r|S_t, a)[r + \gamma V(s')]$$

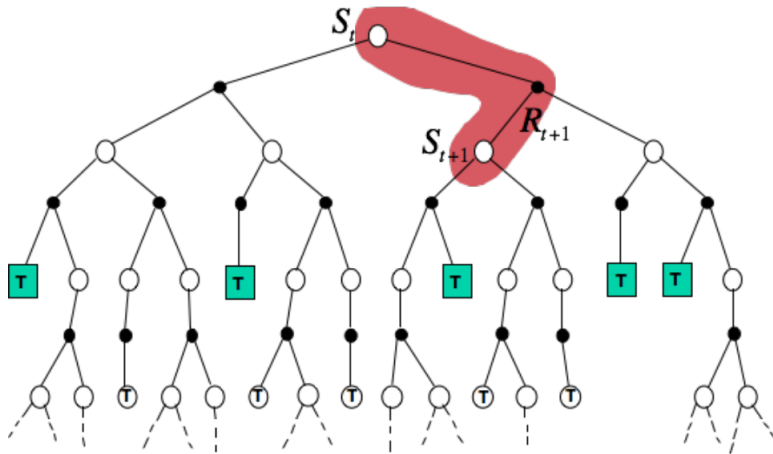


$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



# TD(0) Method

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$





# TD Methods Bootstrap and Sample

- **Bootstrapping:** update involves an estimate
  - MC does not bootstrap
  - DP bootstrap
  - TD bootstrap

# TD Methods Bootstrap and Sample


- **Bootstrapping:** update involves an estimate
  - MC does not bootstrap
  - DP bootstrap
  - TD bootstrap
- **Sampling:** update does not involve an expected value
  - MC samples
  - DP does not sample
  - TD samples

# TD(0) Method

- ▶ **Policy Evaluation** (the prediction problem):
  - for a given policy  $\pi$ , compute the state-value function  $v_\pi$


- ▶ **Remember:** Simple every-visit **Monte Carlo method**:

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

 **target:** the actual return after time  $t$

- ▶ The simplest **Temporal-Difference** method TD(0):

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

 **target:** an estimate of the return

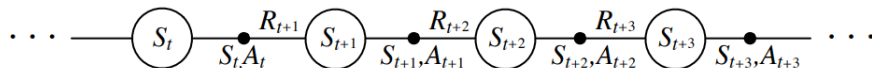
# Table of Contents

1 Recap

2 SARSA

# Learning an Action-Value Function

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



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

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

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

Turn this into a control method by always updating the policy to be **greedy** with respect to the current estimate

## 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

- ▶ Instead of the sample value-of-next-state, use the **expectation!**

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

- ▶ **Expected Sarsa** performs better than Sarsa (but costs more)

- SARSA is an **on-policy** algorithm which means that while learning the optimal policy it uses the current estimate of the optimal policy to generate the behaviour
- SARSA converges to an **optimal policy** as long as all state-action pairs are visited an infinite number of times and the policy converges in the limit to the greedy policy ( $\epsilon = \frac{1}{t}$ ).