# Model-Free
# Temporal-Difference (Q-Learning)

Alina Vereshchaka

CSE4/510 Reinforcement Learning
Fall 2019

*avereshc@buffalo.edu*

September 19, 2019

# Overview

# Table of Contents

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

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

- Incremental every-visit Monte-Carlo:
  - Update value $V(S_t)$ toward actual return $G_t$

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

- Simplest Temporal-Difference learning algorithm: TD(0)
  - Update value $V(S_t)$ toward estimated returns $R_{t+1} + \gamma V(S_{t+1})$

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

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

# 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:

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

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

‣ The simplest Temporal-Difference method TD(0):

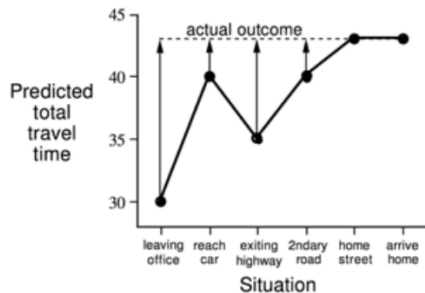$$V(S_t) \leftarrow V(S_t) + \alpha \Big[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \Big]$$

**target**: an estimate of the return

# Example: Driving Home

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|-------|------------------------|----------------------|----------------------|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

Changes recommended by Monte Carlo methods ($\alpha=1$)

Changes recommended by TD methods ($\alpha=1$)

# Table of Contents
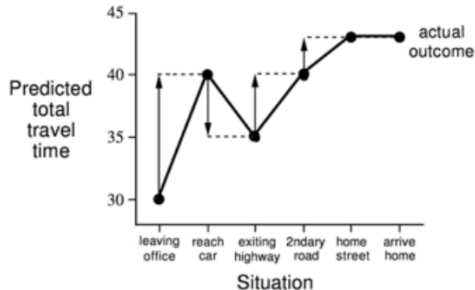
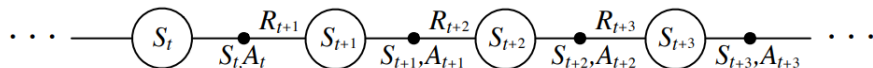▸ Estimate $q_\pi$ for the current policy $\pi$

$$\cdots \quad \underset{S_t, A_t}{\overset{R_{t+1}}{\bullet}} \boxed{S_{t+1}} \underset{S_{t+1}, A_{t+1}}{\overset{R_{t+2}}{\bullet}} \boxed{S_{t+2}} \underset{S_{t+2}, A_{t+2}}{\overset{R_{t+3}}{\bullet}} \boxed{S_{t+3}} \underset{S_{t+3}, A_{t+3}}{\bullet} \quad \cdots$$

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

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

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

# SARSA

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 $Q(terminal, \cdot) = 0$

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)
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma Q(S', A') - Q(S, A) \big]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

---

# SARSA

- 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 ($\epsilon = \frac{1}{t}$).

# Table of Contents

# Q-learning: Off-Policy TD Control

In Q-learning the learned action-value function, Q, directly approximates the optimal action-value function, independent of the policy being followed.

$$Q(s_t, a_t) \leftarrow (s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)$$

# Q-Learning Algorithm

**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 $Q(terminal, \cdot) = 0$

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 $R, S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_a Q(S', a) - Q(S, A) \right]$
        $S \leftarrow S'$
    until $S$ is terminal

Target

Prediction

Immediate Reward

loss