

# **ENGG6600 Reinforcement Learning Assignment-2 Fall 2022**

Yaowen Mei  
1177855

## Problem 1- First-Visit vs Every-Visit Monte Carlo

In the Mars Rover example in the lectures, use  $\gamma = 1$ . Assume the policy is given as below:

- TL in all the states
- S1, and S7 transition to terminal state upon any action.
- 

$$\eta = \underbrace{S_3, TL, 0}_{t_0}, \underbrace{S_3, TL, 0}_{t_1}, \underbrace{S_2, TL, 0}_{t_2}, \underbrace{S_1, TL, 1}_{t_3}, \underbrace{End}_{t_4}$$

	s1	s2	s3	s4	s5	s6	s7
r	1	0	0	0	0	0	10
$\pi$	TL	TL	TL	TL	TL	TL	TL

### Questions List:

1. Use **first-visit** Monte Carlo to estimate the value functions of all the states.
2. Use **every-visit** Monte Carlo to estimate the value functions of all the states.

### SOLUTION-1:

Derive all the Returns for each time stamp ( $i = 1$  for the first iteration):

$$\begin{aligned} G_{i,0}(S_3) &= R_{i,1} + \gamma R_{i,2} + \gamma^2 R_{i,3} + \gamma^3 R_{i,4} = 1 \\ G_{i,1}(S_3) &= R_{i,2} + \gamma R_{i,3} + \gamma^2 R_{i,4} = 1 \\ G_{i,2}(S_2) &= R_{i,3} + \gamma R_{i,4} = 1 \\ G_{i,3}(S_1) &= R_{i,4} = 1 \end{aligned}$$

Along the trajectory, S3, S2, and S1 are visited.

	s1	s2	s3	s4	s5	s6	s7
$\hat{v}(s)$	0	0	0	0	0	0	0
$N(s)$	0	0	0	0	0	0	0
$G(s)$	0	0	0	0	0	0	0

## SOLUTION-1.1 First-Visited MC:

For S1:

$$\begin{aligned} N(S_1) &= N(S_1) + 1 = 0 + 1 = 1 \\ G(S_1) &= G(S_1) + G_{i,3}(S_1) = 0 + 1 = 1 \\ \hat{v}_\pi(S_1) &= G(S_1)/N(S_1) = \frac{1}{1} = 1 \end{aligned}$$

For S2:

$$\begin{aligned} N(S_2) &= N(S_2) + 1 = 0 + 1 = 1 \\ G(S_2) &= G(S_2) + G_{i,2}(S_2) = 0 + 1 = 1 \\ \hat{v}_\pi(S_2) &= G(S_2)/N(S_2) = \frac{1}{1} = 1 \end{aligned}$$

For S3 (S3 were visited twice, but only the first time visit ( $t_0$ ) evokes a training):

$$\begin{aligned} N(S_3) &= N(S_3) + 1 = 0 + 1 = 1 \\ G(S_3) &= G(S_3) + G_{i,0}(S_3) = 0 + 1 = 1 \\ \hat{v}_\pi(S_3) &= G(S_3)/N(S_3) = \frac{1}{1} = 1 \end{aligned}$$

For S4 to S7, they were not learned anything from this trajectory.

	s1	s2	s3	s4	s5	s6	s7
$v^\pi(s)$	1	1	1	0	0	0	0
$N(s)$	1	1	1	0	0	0	0
$G(s)$	1	1	1	0	0	0	0

Table 1: First-visit MC estimation of value functions for all states.

## SOLUTION-1.2 Every-Visit MC:

For S1:

$$\begin{aligned} N(S_1) &= N(S_1) + 1 = 0 + 1 = 1 \\ G(S_1) &= G(S_1) + G_{i,3}(S_1) = 0 + 1 = 1 \\ \hat{v}_\pi(S_1) &= G(S_1)/N(S_1) = \frac{1}{1} = 1 \end{aligned}$$

For S2:

$$\begin{aligned} N(S_2) &= N(S_2) + 1 = 0 + 1 = 1 \\ G(S_2) &= G(S_2) + G_{i,2}(S_2) = 0 + 1 = 1 \\ \hat{v}_\pi(S_2) &= G(S_2)/N(S_2) = \frac{1}{1} = 1 \end{aligned}$$

For S3 (S3 were visited twice, so it will be trained twice:

**The first time visit S3:**

$$\begin{aligned} N(S_3) &= N(S_3) + 1 = 0 + 1 = 1 \\ G(S_3) &= G(S_3) + G_{i,0}(S_3) = 0 + 1 = 1 \\ \hat{v}_\pi(S_3) &= G(S_3)/N(S_3) = \frac{1}{1} = 1 \end{aligned}$$

**The second time visit S3:**

$$\begin{aligned} N(S_3) &= N(S_3) + 1 = 1 + 1 = 2 \\ G(S_3) &= G(S_3) + G_{i,1}(S_3) = 1 + 1 = 2 \\ \hat{v}_\pi(S_3) &= G(S_3)/N(S_3) = \frac{2}{2} = 1 \end{aligned}$$

For S4 to S7, they were not learned anything from this trajectory.

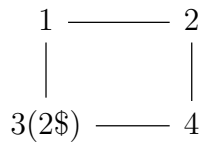
	s1	s2	s3	s4	s5	s6	s7
$v^\pi(s)$	1	1	1	0	0	0	0
$N(s)$	1	1	2	0	0	0	0
$G(s)$	1	1	2	0	0	0	0

Table 2: Every-visit MC estimation of value function for all states.

## Problem 2—Q Learning vs SARSA

Agent A lives in a 2x2 grid as shown below:

The states correspond to the numbered squares. Her possible actions are **MoveNorth**, **MoveSouth**, **MoveEast**, **MoveWest**.



The agent earns \$2 every time she lands in state 3. There are no other rewards or penalties. The reward functions  $R(s, a)$  are given. Note that this is a continuing task

Assuming learning rate  $\alpha = 1$  and discount rate  $\gamma = 0.9$ .

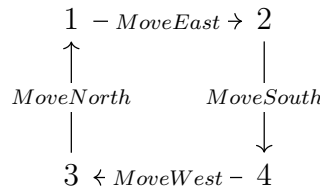
Reward $R(s, a)$	1	2	3	4
MoveNorth	0	0	0	0
MoveSouth	2	0	0	0
MoveEast	0	0	0	0
MoveWest	0	0	0	2

State Value $Q(s, a)$	1	2	3	4
MoveNorth	0	0	0	0
MoveSouth	0	0	0	0
MoveEast	0	0	0	0
MoveWest	0	0	0	0

## Questions List:

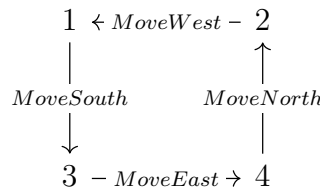
- Agent A starts in square 1 and performs the following actions: MoveEast, MoveSouth, MoveWest, MoveNorth. After each action, the Q-Table is updated using **Q-Learning**, with the the usual update formula:

$$q(s, a) \leftarrow q(s, a) + \alpha \left( r(s, a) + \gamma \max_{a'} q(s', a') - q(s, a) \right)$$



- Agent A **continues** in square 1 and performs the following actions: MoveSouth, MoveEast, MoveNorth, MoveWest. After each of the first three action, the Q-Table is updated using **SARSA**, with the the usual update formula:

$$q(s, a) \leftarrow q(s, a) + \alpha \left( r(s, a) + \gamma q(s', a') - q(s, a) \right)$$



Calculate the Q Table after the last update.

## SOLUTION-2-1 Q-Learning:

- Initial state in Episode 0, First action,
- Current State  $s = 1$ ,

- Current Action  $a = \text{MoveEast}$ ,
- Next State  $s' = 2$
- Update  $q(1, \text{MoveEast})$

$$\begin{aligned} q(1, \text{ME}) &\Leftarrow q(1, \text{ME}) + \alpha \left( r(1, \text{ME}) + \gamma \text{Max} \begin{bmatrix} q(2, \text{MN}) \\ q(2, \text{MS}) \\ q(2, \text{ME}) \\ q(2, \text{MW}) \end{bmatrix} - q(1, \text{ME}) \right) \\ &= 0 + 1 \times (0 + 0.9 \times 0 - 0) = 0 \end{aligned}$$

- Episode 0, Second action
- Current State  $s = 2$ ,
- Current Action  $a = \text{MoveSouth}$ ,
- Next State  $s' = 4$
- Update  $q(2, \text{MoveSouth})$

$$\begin{aligned} q(2, \text{MS}) &\Leftarrow q(2, \text{MS}) + \alpha \left( r(2, \text{MS}) + \gamma \text{Max} \begin{bmatrix} q(4, \text{MN}) \\ q(4, \text{MS}) \\ q(4, \text{ME}) \\ q(4, \text{MW}) \end{bmatrix} - q(2, \text{MS}) \right) \\ &= 0 + 1 \times (0 + 0.9 \times 0 - 0) = 0 \end{aligned}$$

- Episode 0, Third action
- Current State  $s = 4$ ,
- Current Action  $a = \text{MoveWest}$ ,
- Next State  $s' = 3$
- Update  $q(4, \text{MoveWest})$

$$\begin{aligned} q(4, \text{MW}) &\Leftarrow q(4, \text{MW}) + \alpha \left( r(4, \text{MW}) + \gamma \text{Max} \begin{bmatrix} q(3, \text{MN}) \\ q(3, \text{MS}) \\ q(3, \text{ME}) \\ q(3, \text{MW}) \end{bmatrix} - q(4, \text{MW}) \right) \\ &\Leftarrow 0 + 1 \times (2 + 0.9 \times 0 - 0) = 2 \end{aligned}$$

- Episode 0, Forth action (last action)

- Current State  $s = 3$ ,
- Current Action  $a = \text{MoveNorth}$ ,
- Next State  $s' = 1$
- Update  $q(3, \text{MoveNorth})$

$$\begin{aligned}
 q(3, \text{MN}) &\Leftarrow q(3, \text{MN}) + \alpha \left( r(3, \text{MN}) + \gamma \text{Max} \begin{bmatrix} q(1, \text{MN}) \\ q(1, \text{MS}) \\ q(1, \text{ME}) \\ q(1, \text{MW}) \end{bmatrix} - q(3, \text{MN}) \right) \\
 &\Leftarrow 0 + 1 \times (0 + 0.9 \times 0 - 0) = 0
 \end{aligned}$$

Q-Learning Q Table	1	2	3	4
MoveNorth	0	0	0	0
MoveSouth	0	0	0	0
MoveEast	0	0	0	0
MoveWest	0	0	0	2

Table 3: Q-Table after episode-0's last update via Q-Learning method.

## SOLUTION-2-2 SARSA:

- Episode 1, First action under SARSA,
- Current State  $s = 1$ ,
- Current Action  $a = \text{MoveSouth}$ ,
- Next State  $s' = 3$
- Update  $q(1, \text{MoveSouth})$
- Next Action  $a' = \text{MoveEast}$

$$\begin{aligned}
 q(1, \text{MoveSouth}) &\Leftarrow q(1, \text{MoveSouth}) + \alpha (r(1, \text{MoveSouth}) + \gamma q(3, \text{MoveEast}) - q(1, \text{MoveSouth})) \\
 q(1, \text{MoveSouth}) &\Leftarrow 0 + 1 \times (2 + 0.9 \times 0 - 0) = 2
 \end{aligned}$$

- Episode 1, Second action under SARSA,
- Current State  $s = 3$ ,
- Current Action  $a = \text{MoveEast}$ ,

- Next State  $s' = 4$
- Update  $q(3, MoveEast)$
- Next Action  $a' = MoveNorth$

$$q(3, MoveEast) \leftarrow q(3, MoveEast) + \alpha (r(3, MoveEast) + \gamma q(4, MoveNorth) - q(3, MoveEast))$$

$$q(3, MoveEast) \leftarrow 0 + 1 \times (0 + 0.9 \times 0 - 0) = 0$$

- Episode 1, Second action under SARSA,
- Current State  $s = 4$ ,
- Current Action  $a = MoveNorth$ ,
- Next State  $s' = 2$
- Update  $q(4, MoveNorth)$
- Next Action  $a' = MoveWest$

$$q(4, MoveNorth) \leftarrow q(4, MoveNorth) + \alpha (r(4, MoveNorth) + \gamma q(2, MoveWest) - q(4, MoveNorth))$$

$$q(4, MoveNorth) \leftarrow 0 + 1 \times (0 + 0.9 \times 0 - 0) = 0$$

SARSA Q Table	1	2	3	4
MoveNorth	0	0	0	0
MoveSouth	2	0	0	0
MoveEast	0	0	0	0
MoveWest	0	0	0	2

Table 4: Q-Table after episode-1's last update via SARSA.