

Assignment 3

Anushka Bhandari
2016134

Q1 Initialise

$$\pi(s) \in A(s) \quad \text{for all } s \in S$$

$$Q(s, a) \in \mathbb{R} \quad \forall a \in A \text{ and } s \in S$$

$$\text{count}(s, a) = 0 \quad \forall a \in A \text{ and } s \in S$$

~~loop for ever that is for.~~
~~for each episode~~

for each episode : loop for ever

Choose $s_0 \in S$, $A_0 \in A(s_0)$ randomly such that
all pairs have prob > 0

Generate episode from s_0 to A_0 following π

$$\pi : s_0, A_0, R_1, \dots, s_T, A_T, R_{T+1}$$

$$G \leftarrow 0$$

loop for each step $t = T-1, T-2, \dots, 0$:

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

Unless the pair s_t, A_t appear in $s_0, A_0,$
 $s_1, A_1, \dots, s_{t+1}, A_{t+1}$

$$Q(s_t, A_t) = Q(s_{t+1}, A_{t+1}) + \frac{1}{\text{count}(s_t, A_t)} (G - Q(s_t, A_t))$$

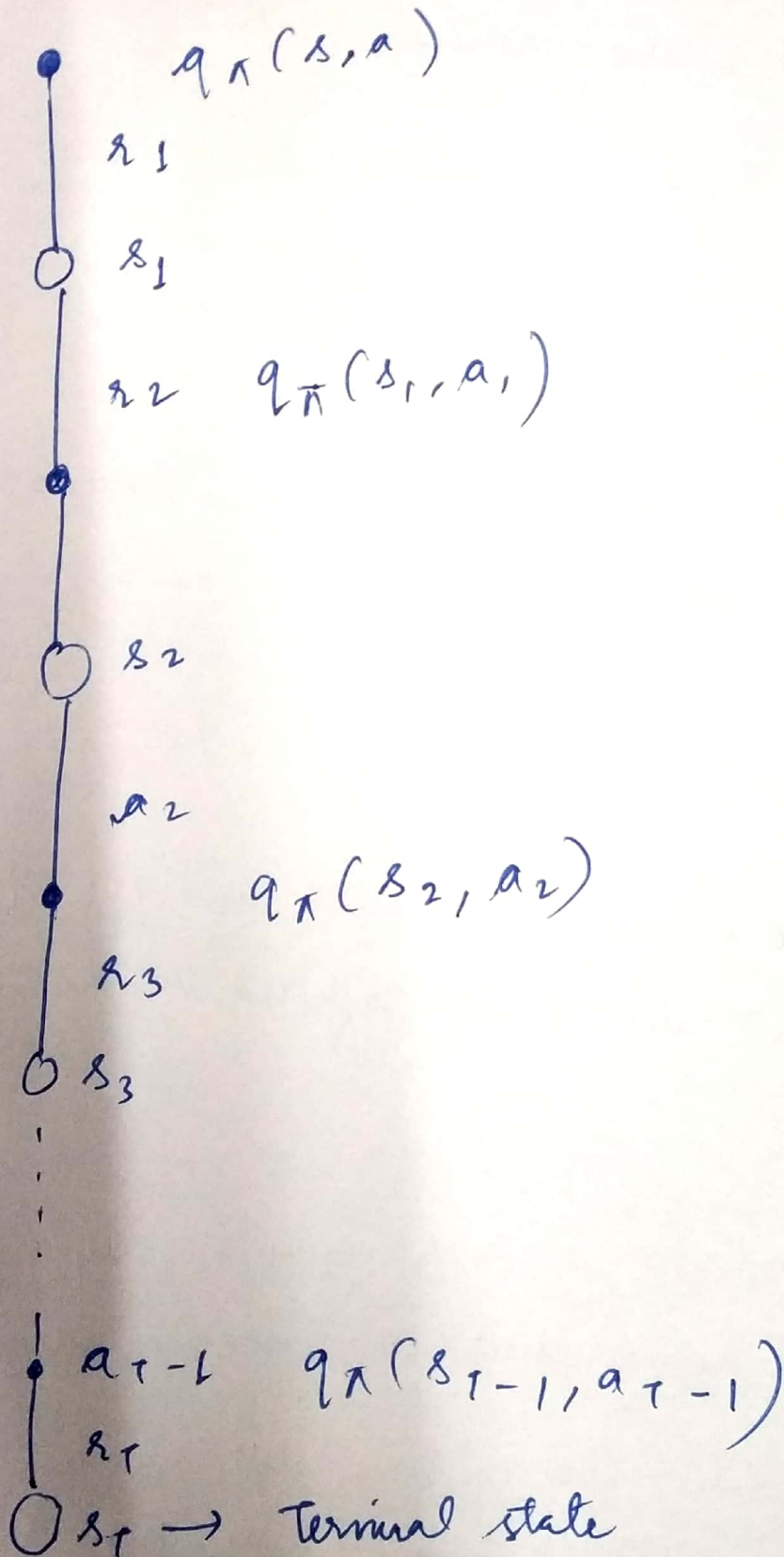
$$\text{count}(s_t, A_t) += 1$$

$$\pi(s_t) \leftarrow \arg \max_a Q(s_t, A_t)$$

Q2

Back up diagram for MC estimation of q_π

• \rightarrow representation action
 ○ \rightarrow representation state



Q3

$$Y(s) = \frac{\sum_{k \in T(s)} \prod_{l: T(l)-1} G_l}{\sum_{k \in T(s)} \prod_{l: T(l)-1}}$$

$T(s, a)$ are set of all ~~timestamp~~ the time steps s and a are visited.
This is for every visit method.

so for
 $Q(s, a)$, $\cdot T(s)$ is replaced

$$Q(s, a) = \frac{\sum_{k \in T(s, a)} \prod_{l+1: T(l)-1} G_l}{\sum_{k \in T(s, a)} \prod_{l+1: T(l)-1}}$$

- Q5
- ① TD learns directly from episodes of experience
 - ② TD is model free. No knowledge of MDP transitions.

lets take a driver has a lot of experience driving from office to home.

OFFICE \longrightarrow lane 1 \longrightarrow lane 2 \longrightarrow HOME

But on the next day suppose lane 2 is closed due to construction, the route then changes to

OFFICE \longrightarrow ^{lane} ~~road~~ 1 \longrightarrow lane 3 \longrightarrow lane 4 \longrightarrow HOME

Had we used MC to update we would have to wait till we reached home.

TD updates at every step.

We can use the learnt values from office and lane 1 from the start

Q6.3

Only the $V(A)$ was changed
This is because the updates that follow the below equation.

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

$\therefore V(s') = 0$ as it is the terminal state

$$R = 0 \text{ and } V(s) = 0.5 \text{ and } \alpha = 0.1$$

$$\begin{aligned} V(A) &= V(A) - \alpha V(A) \\ &= 0.5 - 0.1(0.5) = 0.45 \end{aligned}$$

Only $V(A)$ was updated because the trajectory ended at A and other $V(s)$ balanced out.

Q6.4

Since TD updates using rewards and MC updates using returns, TD will be more affected by α .

Q6.5

Since the $V(s)$ for each s is steadily updated to optimal, ~~there~~ hence there are fluctuations in TD RMS.

After it reaches optimal the updating still continues.

Hence it causes $V(s)$ to move away from $V^*(s)$.

The fluctuation is greater for larger α because the step size derives the value of update.

Q 8

There is only one term difference $\sum_a Q(s', a)$ and $\max_a Q(s', a)$ between SARSA and Q learning

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha (R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t))$$

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha (R_{t+1} + \gamma Q(s', A') - Q(s_t, A_t))$$

In case every selection is greedy with every other condition being the same.

$$\text{then } Q(s', A') = \max_a Q(s', a)$$

This will make the same action selection and hence the same weight update.