# Question 1)

In the step, $Q(S_t, A_t) \leftarrow$ Average $(Returns(S_t, A_t))$ we are taking average every time and thus it involves a lot of calculations.

we can avoid this by just storing the value of $G$ and the number of times the particular state action pair has been encountered.

i.e.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{count[s][a]} \{G - Q(S_t, A_t)\}$$

## Pseudo code:

Initialize:

$\pi(s) \in A(s)$

$Q(s,a) \in \mathbb{R}$

$Count(s,a) = 0 \quad \forall \; s,a$

Loop forever (for each episode) ̶f̶o̶r̶ ̶e̶a̶c̶h̶ ̶e̶p̶i̶s̶o̶d̶e̶ :

̶G̶o̶t̶o̶ Choose $S_0 \in S, \quad A_0 \in A(S_0)$

Generate an episode from $S_0, A_0$, following $\pi$ :

$S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \ldots, 0$ :

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

$Count[S_t][A_t] += 1$

unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, \ldots, S_{t-1}, A_{t-1}$
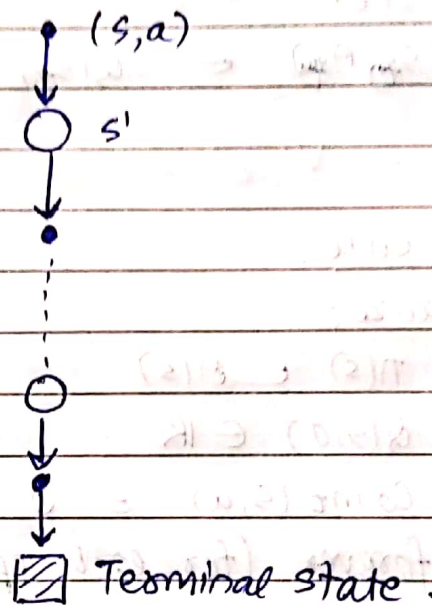
̶A̶p̶p̶e̶n̶d̶ ̶G̶ ̶t̶o̶ ̶R̶e̶t̶u̶r̶

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{count(S_t, A_t)} \times [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \underset{a}{argmax} \; Q(S_t, a)$$

## Question 2:

In Monte Carlo ES, if we use $q(s,a)$ instead of $V(s)$, we generate a state action pair at the start of episode and continue till we reach a terminal state.

This can be shown in the form of backup diagram as follows:



## Question 3

We need to maintain an array of time steps for the pair $(s,a)$, which will store the time instants when that pair is visited for that episode.

Updated formula in terms of $Q(s,a)$ can be written as:

$$Q(s,a) = \frac{\sum\limits_{t \in \tau(s,a)} \gamma^{\tau(t)-1} G_t}{\sum\limits_{t \in \tau(s,a)} \gamma^{\tau(t)-1}}$$

where,

$T(t)$ → first time of termination following time t

$G_t$ → return after t up through $T(t)$

## Question5

TD would work better in an example in which we move to a new building and a new parking lot, which is near the same highway and we have previous experience of driving to the old building. This is because: In this case, only a part of the driving route has changed and some of the states (like highway) still remain the same.

Now, in the TD method, the state values are updated on the fly i.e. without the need to generate the full episode. Thus since the starting values for these states are already close to true values, it will lead to faster convergence.

Yes, same thing will happen in original scenario because of the above reason if initial state values estimate is close to true values.

## Question8:

In case of greedy action selection, Q-Learning will not be the same as SARSA.

This is because the way we update Q in both the methods differs. In Q-Learning the q values are updated irrespective of the action taken and the action is then chosen as per the updated Q value. Whereas, in SARSA next action is chosen according to the current Q value and then Q is updated.

So, since the action chosen might be different in both cases, both the methods perform differently.