Assignment 3.

Reinforcement Learning

① The incremental implementation is given as

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^{n} R_i$$

where $Q_{n+1}$ = estimate of its action value after it has been selected $n$ times.

$$\therefore \quad Q_{n+1} = Q_n + \frac{1}{n}[R_n - Q_n] \quad \Rightarrow \quad R_n = n^{th} \text{ reward.}$$

The pseudo code given in the book is in efficient because, for each state-action pair, it maintains a list of all returns and repeatedly calculates their mean.

So the efficient code can be written using the concept of incremental implementation.

Also we can write that

$$Q_n(S_t, A_t) = \frac{1}{n} \sum_{i=n}^{n} G_i(S_t, A_t)$$

$$= Q_{n-1}(S_t, A_t) + \frac{1}{n}[G_n(S_t, A_t) - Q_{n-1}(S_t, A_t)]$$

## Pseudo code :-

Initialize :

$\pi(s) \in A(s) \quad \forall s \in S$

$Q(s,a) \in R \qquad \forall s \in S, a \in A(s)$

Returns $(s,a) \leftarrow$ empty list $\forall s \in S, a \in A(s)$

$n(s,a) \leftarrow 0 \quad \forall s \in S, a \in A(s)$

Loop forever :

Choose $S_0 \in S, A_0 \in A(S_0)$ randomly such that all pairs have probability $> 0$.

Generate an episode from $S_0, A_0$. following $\pi$: $S_0, A_0, R_1 \cdots$

$\qquad S_{T-1}, A_{T-1}, R_T$.

$G \leftarrow 0$

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

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

Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1 \cdots S_{t-1}, A_{t-1}$

$\qquad n(S_t, A_t) = n(S_t, A_t) + 1.$

$\qquad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \dfrac{1}{n(S_t, A_t)} \left[ G - Q(S_t, A_t) \right] \qquad$ ——①

$$\pi(s_t) \leftarrow \underset{a}{\arg\max} \, Q(s_t, a)$$

OR

Instead of using ①, we can write in terms of rewards.

ie

$$Returns(s_t, A_t) \leftarrow Returns(s_t, A_t) + \frac{(G - Returns(s_t, A_t))}{n(s_t, A_t)}$$

$$Q(s_t, A_t) \leftarrow Returns(s_t, A_t).$$

$$\pi(s_t) \leftarrow \underset{a}{\arg\max} \, Q(s_t, a)$$

In both the cases, Monte Carlo ES is satisfied because the policy improvement is acheived at all the states visited in the episodes.

Instead of using more space to store all the return values to find the mean, an alternate formula is used to find the optimal policy $\pi$. using incremental implementation.
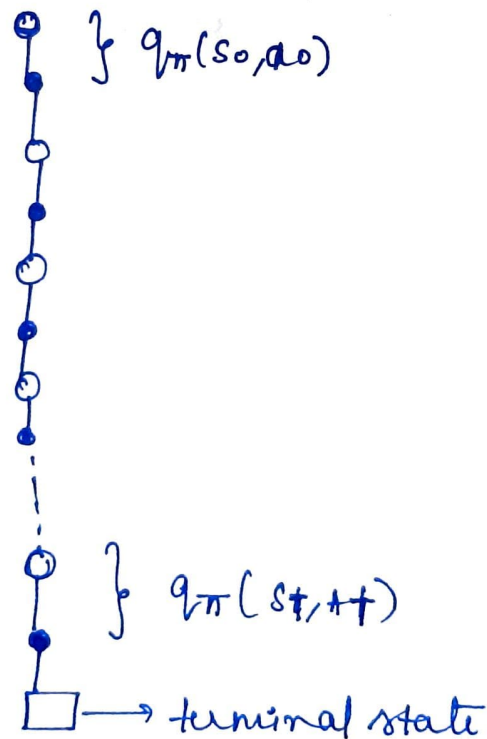
② Back up diagram is the graphical representation of algorithm by representing state, action, state transition, reward etc. Value function is transferred back to a state from its successor state or state action.

$\circ \longrightarrow$ State value

$\bullet \longrightarrow$ State-action value

$\begin{smallmatrix}\circ\\|\end{smallmatrix} \longrightarrow$ Action.

.. the backup diagram for monte carlo for $q_\pi$ is.

$\begin{smallmatrix}\circ\\|\\\bullet\end{smallmatrix} \} q_\pi(s_0, a_0)$

$q_\pi(s_t, a_t)$

$\square \longrightarrow$ terminal state

③ Given a starting state $s_t$, the probability of the subsequent state-action trajectory, $A_t, S_{t+1}, A_{t+1}, \dots S_T$, occuring under any policy $\pi$ is

$$\Pr\{A_t, S_{t+1}, A_{t+1}, \dots S_T | s_t, A_{t:T-1} \sim \pi\}$$

$$= \pi(A_t|s_t) \, p(S_{t+1}|s_t, A_t) \, \pi(A_{t+1}|s_{t+1}) \dots p(S_T|S_{T-1}, A_{T-1})$$

$$= \prod_{k=t}^{T-1} \pi(A_k|s_k) \, p(S_{k+1}|s_k, A_k)$$

where $p$ is the state-transition probability function. The importance sampling ratio is

$$\rho_{t:T-1} \doteq \frac{\prod_{k=t}^{T-1} \pi(A_k|s_k) \, p(S_{k+1}|s_k, A_k)}{\prod_{k=t}^{T-1} b(A_k|s_k) \, p(S_{k+1}|s_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|s_k)}{b(A_k|s_k)}$$

$b$ = Behavioural policy $\pi$ = Target policy.

The weighted importance sampling, which uses a weighted average is

$$V(s) = \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} \, G_t}{\sum_{t \in \mathcal{T}(s)} \rho_{t:T-1}}$$

The state-action value $Q(s,a)$ can be written as

$$Q(s,a) = \sum_{t \in T(s,a)} R_t + \frac{\sum_{t \in T(s)} t_{t+1:T(t)-1} G_{t+1}}{\sum_{t \in T(s)} t_{t+1:T(t)-1}}$$

where $T(s,a)$ is the state action pair of all time steps.

## Blackjack game.

fig 5.1 → Code ⟹ $Qs\_4\_fig5-1.m$

fig 5.2 → code ⟹ $Qs\_4\_fig5-2.m$.

⑤ Exercise 6.2.

In the scenario discussed, for the first drive, TD will be more efficient. TD will start updating the states immediately when the person starts driving. Each updated state accelerate the one prior to it. But in Monte carlo, the updates are done only after the journey. There will not be immediate update. MC bring an average improvement of all the states which are not even needed.

After the first drive, TD will use the past estimates to update the new states which will be efficient in this type of scenario.

⑥ Code: Qs_6_1.m & Qs_6_2.m

Figure: Qs_6_left.png & Qs_6_right.png

6.3) In the figure, for the first episode, the state is at E. The value update for the state is that all states are initialized to ~~zero~~ 0.5.

$$V(A) = V(A) + \alpha[R - 0.5] \quad ; \quad 0.5 + 0.1[0 - 0.5]$$

$$\circ \underline{0.45}$$

As reward = 0 & gamma = 1, the state changes its values to non-zero at A or E.

6.4) Using a large value of $\alpha$ will not effect the results. Small alpha values are better for convergence in TD & MC methods. Using higher values of alpha leads to noise ie the previous outcomes changes the values of the states. There occur an error which is undesirable.

But for small values of alpha, the algorithm converge slowly but in an efficient manner ie without noise.

As they emerge slowly, it is not that small $\alpha$ is better that high alpha, But it produces an efficient in a long run and the error will be very less.

6.5) Yes, I think it is because of the initialization of approximate value functions. The identical initialization of all the states, causes a sudden shoot at many states close to each other at the same time. If same initialization is done, we have to simultaneously check for estimates more than some states and less than some other states. This will cause a sudden decrease, and RMS error rise eventually of RMS error.

⑧ Even for greedy policy, Q-learning & SARSA algorithms will be different. The main difference is in SARSA, the ~~optimal~~ greedy action is used to update the state-action value of that state. But in Q-learning, the updated state action value of that state is considered to rechoose the greedy action for the next state.

In action selection, Q learning is off-policy and use any behavioural policy for action selection. But in SARSA is on-policy and uses greedy action selections where exploration is not done.

④ Black jack game.

Figure 5.1 :—

    Code :— as_4_fig5-1. m.

    figure :— as_4_fig5-1.png
             as_4_fig5-1-1-png
             as_4_fig5-1-2.png
             as_4_fig5-1-3.png

Figure 5.2 :—

    Code :— as_4_fig5-2. m.

    figure :— as_4_fig5-2.png.

             as_4_fig5-2-1.png.