

Homework 3

 Q

Ans 1

Let return count (S_t, A_t) denotes the number of returns of (S_t, A_t) we have got so far.

$$Q(S_t, A_t) = \frac{\sum_{r \in \text{Returns}(S_t, A_t)} r}{\text{count}(S_t, A_t)}$$

Let r' be the last return of

Now suppose we get a new return r' , then let the new Q be $Q'(S_t, A_t)$

$$Q'(S_t, A_t) = \frac{\sum_{r \in \text{Returns}(S_t, A_t)} r + r'}{\text{count}(S_t, A_t) + 1}$$

$$Q'(S_t, A_t) = \frac{Q(S_t, A_t) \text{count}(S_t, A_t) + r'}{\text{count}(S_t, A_t) + 1}$$

and we will update $\text{count}(S_t, A_t)$ by 1

∴ pseudo code would be →

Initialize :

$\pi(s) \in A(s)$ (arbitrarily), for $\forall s \in S$

$Q(s, a) \cancel{=} 0$, $\forall s \in S, a \in A(s)$

$\text{count}(s, a) = 0$

Loop forever (for each episode) :

choose $s \in S, a \in A(s_0)$ randomly

generate episode from s_0, a_0, \dots

$G \leftarrow 0$

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

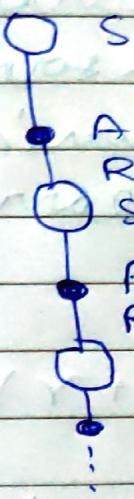
$$Q = QY + R_{t+1}$$

Update the pair S_t, A_t appearance $S_0 A_0, S_1 A_1, \dots, S_{T-1} A_{T-1}$

$$Q(S_t, A_t) = \frac{Q(S_t, A_t) \text{count}(S_t, A_t) + Q}{\text{count}(S_t, A_t) + 1}$$

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

$$\pi(S_t) = \arg\max_a Q(S_t, a)$$



Aus

TD learning will be better than MC learning if the predicted time to go is very close to the actual value.

Now suppose you have very good estimates (almost true) of predicted time to go for all the states.

Now you change your office, in this case only the first state & 2nd state would change but the estimation of predicted time from the moment he enters highway and thereafter will remain same.

Thus in this case TD learning would be better than Monte Carlo learning.

The update in TD learning is →

$$V(S+) = V(S+) + \alpha [R_{++1} + \gamma V(S_{++1}) - V(S+)]$$

Here better is the estimation of state $V(S_{++1})$, better would be estimation of $V(S+)$.

Thus if we know the true value of $V(S_{++1})$ we can easily calculate the true value of $V(S+)$. Thus, the ~~example~~ situation we considered is better for TD.

Aus

Q No, Q-learning and SARSA wouldn't be same in case of Greedy selection.

This is because according to the pseudo code of both the methods, SARSA first chooses the action and then updates whereas Q-learning first updates the value

of or another chooses the action.