

HW-3

1) Initialize

$\pi(s) \in A(s)$ (arbitrarily) for all $s \in S$
 $Q(s, a) \in \mathbb{R}$ (arbitrarily) for all $s \in S, a \in A(s)$
 $\text{Count}(s, a) \in \text{Integers}$ for all $s \in S, a \in A(s)$
 $Q(s, a) \leftarrow 0, \text{Count}(s, a) \leftarrow 0 \quad \forall s \in S, a \in A(s)$

Loop Forever (for each episode):

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

Generate an episode from S_0, A_0 , following π :

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

$G \leftarrow 0$

For $t = T-1, T-2, \dots, 0$

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

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

$A_1, \dots, S_{t-1}, A_{t-1}$

$Q(S_t, A_t) \leftarrow \frac{Q(S_t, A_t) \cdot \text{Count}(S_t, A_t) + G}{\text{Count}(S_t, A_t) + 1}$

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

This pseudocode is equivalent as we have assumed that $Q(S_t, A_t)$ already has the average ~~returns~~ returns for state S_t and A_t for t , say n points

$$\Rightarrow Q(S_t, A_t) = \frac{G_1 + G_2 + \dots + G_n}{n}$$

Now, to update $Q(S_t, A_t)$ for a new point G_{n+1} , we get $Q(S_t, A_t)^u$

$$Q(S_t, A_t)^u = \frac{G_1 + G_2 + \dots + G_n + G_{n+1}}{n+1}$$

Here, we maintained a count in the count (S, A) which was n

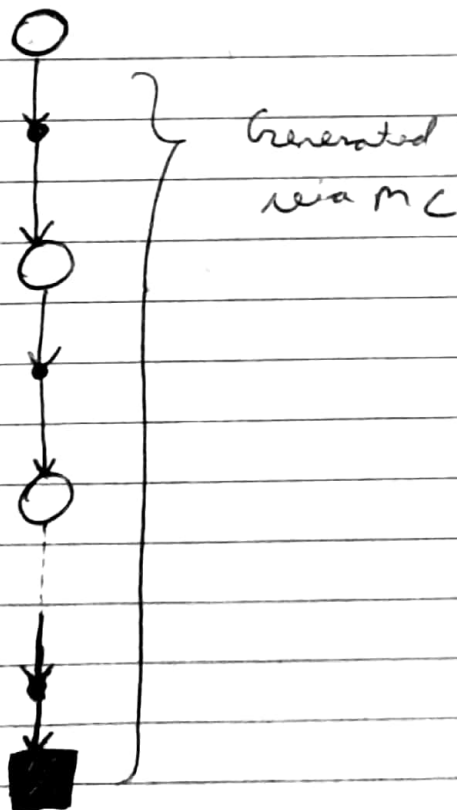
$$\Rightarrow Q(S, A) = \frac{n \cdot Q(S, A) + G_{n+1}}{n+1}$$

$$= \frac{\text{Count}(S, A) \cdot Q(S, A) + G_{n+1}}{\text{Count}(S, A) + 1}$$

\Rightarrow We have maintained the average by devising a recursive relation

\therefore It is the same as mentioned pseudocode

2)



$$3 > P[\cancel{A_t}, S_t, A_t, S_{t+1}, A_{t+1}, S_{t+2}, A_{t+2}, \dots, S_T \mid S_t = s_t, A_t = a_t, A_{t+1} : T-1 \sim \pi]$$

$$= P[S_t | A_t] \cdot P[S_{t+1} | S_t, A_t] \pi(A_{t+2} | S_{t+2}) \dots p[S_T | S_{T-1}, A_{T-1}]$$

$$= P[S_{t+1} | S_t, A_t] \cdot \pi(A_{t+2} | S_{t+2}) \dots P[S_T | S_{T-1}, A_{T-1}]$$

$$= \left(\prod_{k=t+1}^{T-1} \pi(a_k | s_k) p(s_{k+1} | s_k, a_k) \right) p(s_{t+1} | s_t, a_t)$$

$$\Rightarrow \text{Ratio} = \left(\prod_{k=t+1}^{T-1} \pi(a_k | s_k) p(\cancel{s_{k+1}} | s_k, a_k) \right) p(\cancel{s_{t+1}} | s_t, a_t)$$

$$\left(\prod_{k=t+1}^{T-1} h(a_k | s_k) p(s_{k+1} | s_k, a_k) \right) p(s_t | s_{t-1}, a_{t-1})$$

$$= \prod_{k=t+1}^{T-1} \pi(a_k | s_k)$$

$$= \prod_{k=t+1}^{T-1} h(a_k | s_k)$$

$$\Rightarrow Q(s, a) = \frac{\sum_{t \in \tau(s, a)} \prod_{t+1:T-1} G_t}{\sum_{t \in \tau(s, a)} \prod_{t+1:T-1} 1}$$

5) In a nutshell, MC estimates need a lot of episodes in order to make accurate estimates about a new state introduced in a system with previous experience and different set of states while TD uses previous experience to give initial set of estimates with low variance.

$$\text{MC Update: } V_\pi(s_t) = V_\pi(s_t) + \alpha [G_t - V_\pi(s_t)]$$

For the new new states using the new reward and previous state estimates

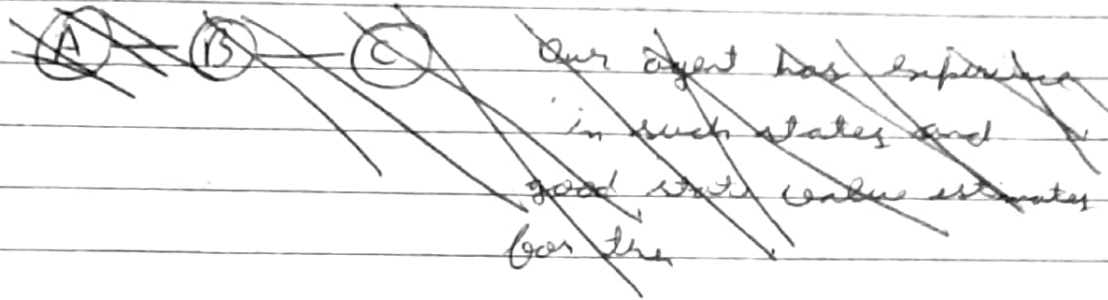
Now, G_t may have high variance

$$\text{TD Update: } V_\pi(s_t) = V_\pi(s_t) + \alpha [R_{t+1} + \gamma V_\pi(s_{t+1}) - V_\pi(s_t)]$$

Here, $V_\pi(s_{t+1})$ has low variance since it has a good estimate of it \rightarrow our assumption

Another case in which TD helps is when it is hard to reach the terminal state. In such cases MC becomes infeasible but TD can still learn significantly.

~~Example:~~



Example:

For the highway instance, these points hold true as with the introduction of the new building, the initial estimates of MC would have extremely high variance due to lack of experience in that state while TD, on the other hand, would use previous experience and current reward to give a low variance estimate for the new building which would serve as great initial estimates.

- 8) The would not be the same in some cases. Reason being is the order of updates and action selection done in both the algorithms. In SARSA, we pick our next action A^* and then update our current state-action ~~pair~~ pair. On the other hand, in Q-learning, we first make updates to the Q-value of current state action pair and then greedily choose from the updated set of Q-values, say that action is $A^{*'}.$ So, it is definitely not guaranteed that $A^* = A^{*'},$ hence the algorithms are different.

6) ~~6.1~~ 6.1) We can conclude that the episode ended at the left-most terminal. Say we have the following MRP,

$C, 0, S_1, 0, S_2, 0, \dots, S_n, 0, A, 0$

~~Initial estimates for all states are 0.5~~
 ~~$u_\pi(S_i) \leftarrow u_\pi(S_i) + \alpha (0 + u_\pi(S_i) - u_\pi(S_i))$~~
~~i.e no change~~

$$u_\pi(C) \leftarrow u_\pi(C) + \alpha (u_\pi(S_1) + 0 - u_\pi(C))$$

~~$u_\pi(A) \leftarrow u_\pi(A)$~~ & states, the initial value is equal

$$\therefore u_\pi(C) \leftarrow u_\pi(C) \quad \text{i.e no change}$$

For, S_i to S_{n-1}

$$u_\pi(S_i) \leftarrow u_\pi(S_i) + 0.1 (0 + u_\pi(S_{i+1}) - u_\pi(S_i))$$

$$\Rightarrow u_\pi(S_i) \leftarrow u_\pi(S_i) \quad \text{i.e no change}$$

For S_n

$$u_\pi(S_n) \leftarrow u_\pi(S_n) + 0.1 (0 + u_\pi(A) - u_\pi(S_n))$$

$$\Rightarrow u_\pi(S_n) \leftarrow u_\pi(S_n) \quad \text{i.e no change}$$

$A \rightarrow A$

$$u_\pi(A) \leftarrow u_\pi(A) + 0.1 (0 + 0 - u_\pi(A))$$

$$\Rightarrow u_\pi(A) \leftarrow u_\pi(A) \quad \text{i.e no change}$$

c) $\alpha_x(A)$ value changes by $\frac{1}{10} \times 0.5 = 0.05$

<6.4> The α ranges are sufficient as for both TD and MC, we have obtained results which converge to the best solution possible for the methods. For MC, the curves get noisy with ~~some~~ increase in α and hence 0.04 is a sufficient place to stop. For $\alpha = 0.05$ in case of TD, we see that the system converges to a sufficiently low ~~absolute~~ error over a period of episodes, hence α ranges are sufficient to characterize the analysis and there is no such ~~extreme difference~~ for alpha that can give a significantly different result than these

<6.5> The error increases because the estimate of $\alpha_x(C)$ changes from its initial estimate in between, which was the ground truth value of $\alpha_x(C)$. So, initially, when there were less episodes, then $\alpha_x(C)$ might have been stable near 0.5, thus causing a decrease in error but as episodes moved on, the value of C may have started to change significantly, thus causing an increase in error from that part. This phenomenon happens, clearly ~~not~~ due to the initialization of state values, especially $\alpha_x(C)$.