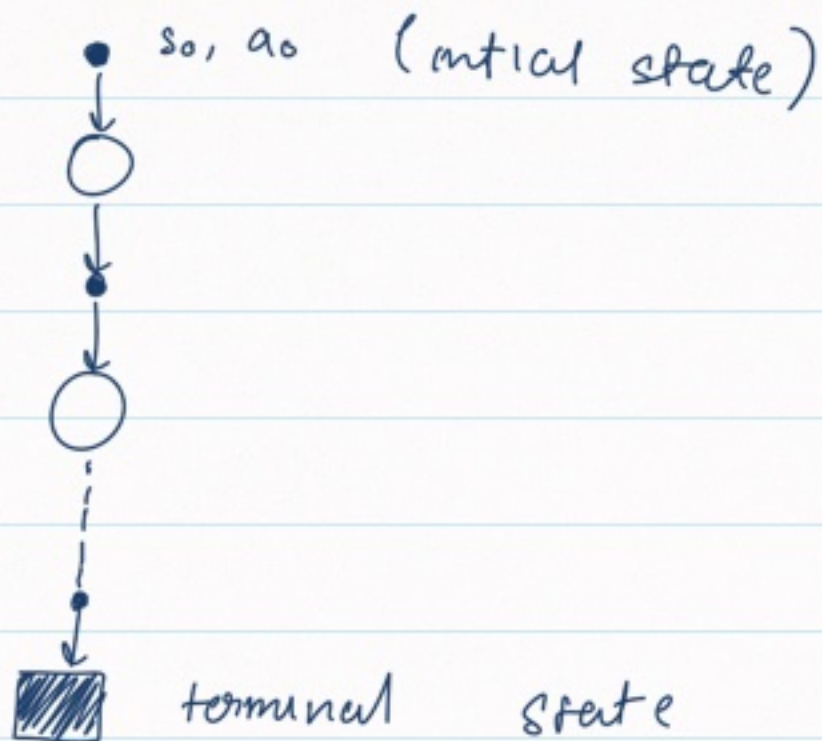


RL HW-3 Theory

Name :- Vishal Samal

Rollno :- 2017209.

2.



5. For the given example where the path is now different initially TD will perform better than MC. This is because the update rule for MC is,

$$V(s_t) \leftarrow V(s_t) + \alpha [G_t - V(s_t)]$$

which would require us to have an entire sequence to calculate state values,

whereas, for TD the update rule is,

$$V(s_t) \leftarrow V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

which requires only the step reward to do this.

Due to this initially TD will be better.

8. Update rule for SARSA,

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma Q(s', a') - Q(s, a)]$$

where $a \in A(s)$, $a' \in A(s')$

and they selected using ϵ -greedy,

Update rule for Q-Learning,

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma \max_a (Q(s', a)) - Q(s, a)]$$

where $a \in A(s)$

and selected using ϵ -greedy.

Now, if ϵ -greedy has $\epsilon = 0$
that is, we are following greedy
policy. The all action values are
optimal action values,

$$\Rightarrow Q(s', a') = \max_a (Q(s', a))$$

Thus, they both will be equal.

P.T.O

1. Pseudo code,

Initialise,

$$\pi(s) \in A(s) \quad [\text{arbitrary}] \quad \forall s \in S$$

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

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

Loop forever (for each episode):

choose $s_0 \in S$, $A_0 \in A(s_0)$ randomly

Generate episode from s_0, A_0

$$G_t \leftarrow 0$$

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

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

$$N(s_t, A_t) \leftarrow N(s_t, A_t) + 1$$

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{1}{N(s_t, A_t)} [G_t - Q(s_t, A_t)]$$

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

We would a variable $N(s_t, a)$

which stores how many times we visit, any state-action pair.

From section 2.4 we know the update rule and we can incorporate that for MC by replacing rewards with returns.

3. For given target policy $\pi(a|s)$ and behaviour $b(a|s)$,

$$V_b(s) = E_b [G_{t+1} | s_t = s]$$

If importance sampling ratio is,

$$P_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

$$\text{Then, } V_{\pi}(s) = E_{\pi} [P_{t:T-1} G_{t+1} | s_t = s]$$

Similarly,

$$Q_{\pi}(s, a) = E_{\pi} [P_{t:T-1} G_{t+1} | s_t = s, A_t = a]$$

We can define a set of all time steps where s, a action pair was visited, $\mathcal{T}(s, a)$

then, ordinary importance sampling,

$$Q(s, a) = \frac{\sum_{t \in \mathcal{T}(s, a)} P_{t:T-1} G_t}{|\mathcal{T}(s, a)|}$$

And, weighted importance sampling,

$$Q(s, a) = \frac{\sum_{t \in \mathcal{T}(s, a)} P_{t:T-1} G_t}{\sum_{t \in \mathcal{T}(s, a)} P_{t:T-1}}$$