

RL HW 3

Vishal Vdandaraa

2016/19.

1) Ex. 5.4

Assume that we track the returns for each state-action pair (s, a) by G_i .

For updating the action value function, we need to average out over G_i 's.

Considering first-visit MC \Rightarrow

$$Q_{n+1}(s, a) = \frac{1}{n} \left(\sum_{i=1}^n G_i \right)$$

where Q_1 is initialized arbitrarily for every state-action pair.

Now, we know \Rightarrow

$$\begin{aligned} Q_{n+1}(s, a) &= \frac{1}{n} \left(\sum_{i=1}^{n-1} G_i + G_n \right) \\ &= \frac{1}{n} \left(\frac{n-1}{n-1} \sum_{i=1}^{n-1} G_i + G_n \right) \\ &= \frac{1}{n} \left((n-1) \cdot Q_n(s, a) + G_n \right) \end{aligned}$$

$$\Rightarrow Q_{n+1}(s, a) = Q_n(s, a) + \frac{1}{n} (G_n - Q_n(s, a))$$

We clearly see from this modified update rule that the count of number of times the particular state-action pair has been encountered and the previous action value function suffice to form the updated action

The updated pseudo code will be \Rightarrow

Initialize:

$\pi(s) \in A(s)$ (arbitrarily) for all $s \in S$.

$Q(s, a) \in \mathbb{R}$ (arbitrarily) for all $s \in S, a \in A(s)$

$\text{Returns}(s, a) = 0$ for all $s \in S, a \in A(s)$

Loop Forever: (for each episode)

Choose $s_0 \in S, A_0 \in A(s_0)$ randomly s.t. all pairs have $^{\text{non}}$ 0 probability.

Generate an episode from s_0, A_0 following

$\pi: s_0, A_0, R_1, \dots, s_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

$\text{Returns}(s, a) + 1$

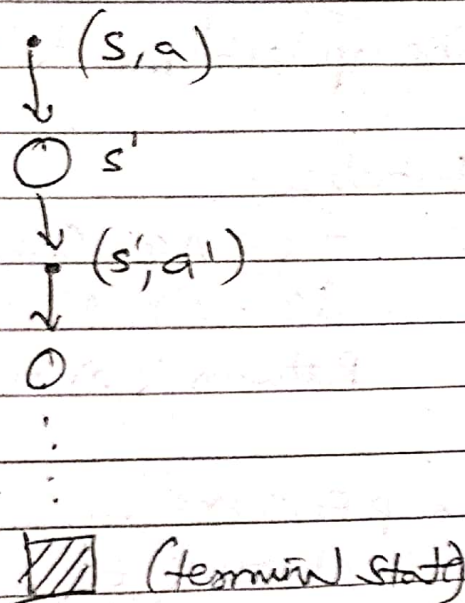
Unless pair s, a appears in $s_0, A_0, \dots, s_{t-1}, A_{t-1}$

$Q(s, a) \leftarrow Q(s, a) + \frac{1}{\text{Returns}(s, a)} (G - Q(s, a))$

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

The above pseudo code is the same as the MC & control ES pseudo code as we have previously shown that both the update rules are equivalent.

2) The backup diagram for q_T is \Rightarrow



The top most node shows the initial state-action pair of a particular episode, the rest of the nodes are state nodes and state-action pairs based on the episode.

Black dots depict state-action pair nodes
White circles represent state nodes.

3) Assume the set of all time-steps in which the state-action pair (s,a) is visited to be denoted as $\Rightarrow T(s,a)$

Then, we can represent $Q(s,a)$ as \Rightarrow

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

Where G_t is the return at steps t to $T(t)$.

and $p_{t:T(t)-1}$ are the importance sampling ratios.

5) In this particular case TD updates would perform better due to \Rightarrow

i) As TD bootstraps, the original state values are retained.

ii) As only initial route in the problem has changed, there will be many common states in the new problem and old problem.

iii) State-value function of TD for original problem would have converged to optimal values since there was a lot of prior experience while driving home from work.

ix) These state value function can be used in the new problem as they would be ~~a~~ very close to ~~the~~ the true state values of the new problem. Hence, There will be faster convergence to true state values.

Even, in the original scenario if the initial value estimates are ~~a~~ close to true values, then these would be faster convergence.

2) If action-selection procedure is greedy, Q -learning is not the same as the sarsa algo. This is because sarsa picks the next action from a given ~~state-action~~ next state according to the Q -values. The update to the Q -values is then based on this next action.

However, in Q -learning, first the Q -values are updated and then the next action is picked according to the updated Q -value. This is why the Q -learning is different from sarsa.

c) Exercise 6.3

We are given the following hyperparameters \Rightarrow

$$\alpha = 0.1$$

$$\gamma = 1 \text{ (undiscounted)}$$

We know, update to state-value is \Rightarrow

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

Since all rewards are 0 except for transitions to terminal state \Rightarrow

$$V(s_t) = V(s_t) + 0.1 (0 + V(s_{t+1}) - V(s_t))$$

Now, all $V(s_t)$ initially were 0.5

$$\Rightarrow V(s_t) = V(s_t) + 0.1 (0.5 - 0.5)$$

$$\Rightarrow V(s_t) = V(s_t)$$

For terminal states \Rightarrow

$$\boxed{V(s_t) = 0}$$

In the first episode, when we go to A that is the left terminal state, our reward is 0.

$$\Rightarrow V(A) = V(A) + 0.1(0 + 0 - V(A))$$

$$\begin{aligned}\Rightarrow V(A) &= V(A) - 0.1 V(A) \\ &= 0.9 V(A) \\ &= 0.9 \times 0.5 \\ &= \underline{\underline{0.45}}\end{aligned}$$

Therefore, we conclude that only state value of A changed in first episode i.e. it decreased to 0.45.

Therefore, the amount of change was 0.1 times the original value = 0.05

Exercise 6.4

Yes, the conclusions about the better performing algo will be affected if we use ~~now~~ a wider range of alpha values. Suppose we take large value of ~~alpha~~ alpha, then even TD(0) will start performing poorly since oscillations are likely to happen as the step size α being large causes larger jumps and hence it will be hard for the convergence to be achieved.

Similar is the case with MC algo. If we pick a wider set of alphas

Yes, there will be a specific α at which the algo might perform better. This can be found out by doing an α -tuning to find out the best α at which RMSE is reduced.

Exercise 6.5

For high α values, there will be oscillations in the convergence (~~with~~ which can be seen as a gradient descent update rule). Therefore, for high ~~oscillations~~ α values, when convergence is achieved, oscillations will occur.

No, this is not a function of how the value function estimates are initialised. At all high α values, oscillations will occur.