

# **Reinforcement Learning – CSE564**

## Assignment 3

Suchet Aggarwal  
2018105

## 1. Pseudo Code:

Initialize:

$\pi(s)$  in  $A(s)$  (arbitrarily), for all  $s$  in  $S$   
 $Q(s, a)$  in  $R$  (arbitrarily), for all  $s$  in  $S$ ,  $a$  in  $A(s)$   
 $\text{counts}(s, a) = 0$ , 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 such that all pairs have probability  $> 0$   
 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,  $t = T-1, T-2, \dots, 0$ :

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

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

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

$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + (G - Q(S_t, A_t)) / \text{counts}(s, a)$

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

Here count stores the number of times the pair of state action has been considered for every state  $s$  and every action  $a$ .

Explanation:

From Chapter 2.4:

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

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

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

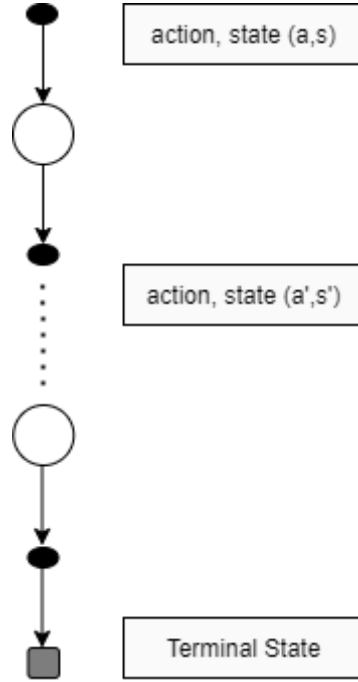
$$Q_{n+1} = \frac{R_n}{n} + \frac{(n-1) \sum_{i=1}^{n-1} R_i}{(n-1) * n}$$

$$Q_{n+1} = \frac{R_n}{n} + \frac{(n-1)Q_n}{n}$$

$$Q_{n+1} = Q_n + \frac{1}{n}(R_n - Q_n)$$

Where  $Q_n$  is the mean value for  $n$  steps and  $n$  is the number of steps

2.



The Backup diagram for Monte Carlo would remain same for Q, as the diagram for V. As for computing the value  $Q(s, a)$  we compute it by generating an episode, and then update the value for pair  $s, a$  using the transitions within that episode.

3. For Q, the first action is fixed for a given state action pair. Thus the important sampling factor becomes:

The probability of state action trajectory under the policy  $\pi$ , starting in state  $S_t$  and taking action  $A_t$

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

Therefore, the important sampling factor becomes:

$$factor_{t:T-1} = \frac{p(S_{t+1} | S_t, A_t) \prod_{k=t+1}^{T-1} p(S_{k+1} | S_k, A_k) \pi(A_k | S_k)}{p(S_{t+1} | S_t, A_t) \prod_{k=t+1}^{T-1} p(S_{k+1} | S_k, A_k) b(A_k | S_k)}$$

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

$$factor_{t:T-1} = \rho_{t+1:T-1}$$

Thus the Q estimate becomes:

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

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

4. Attached in ipynb.

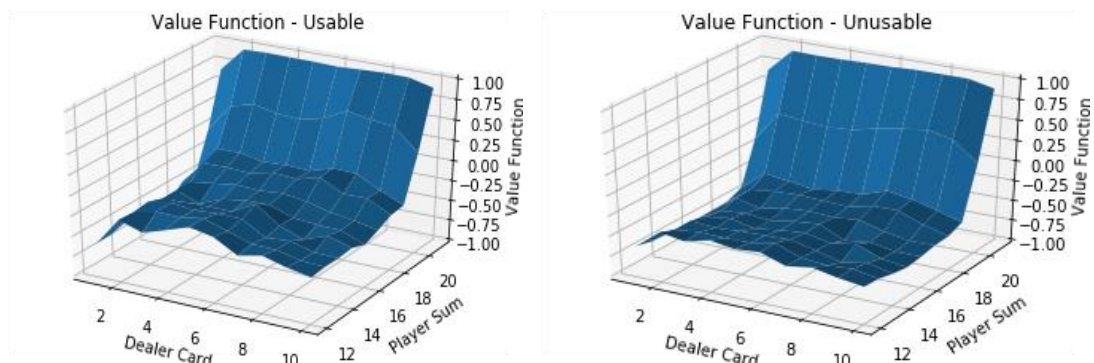


Fig 5.1

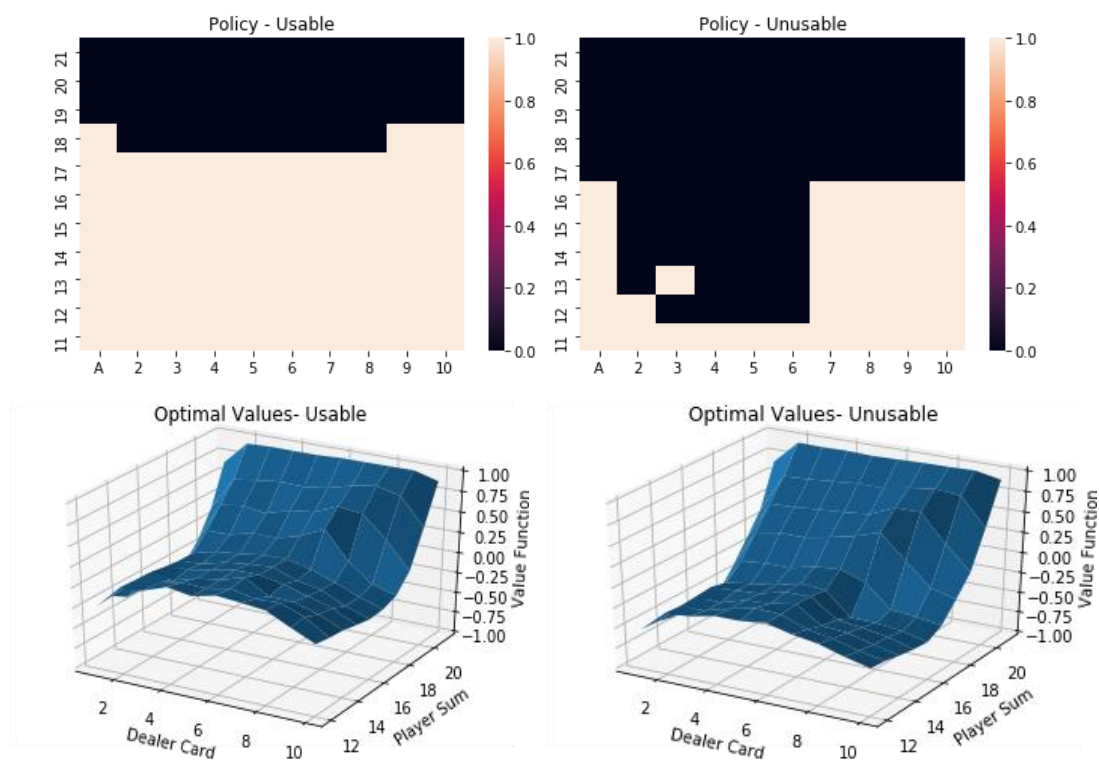


Fig 5.2

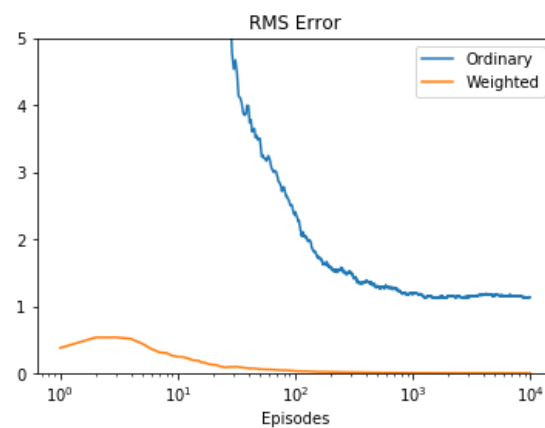


Fig 5.3

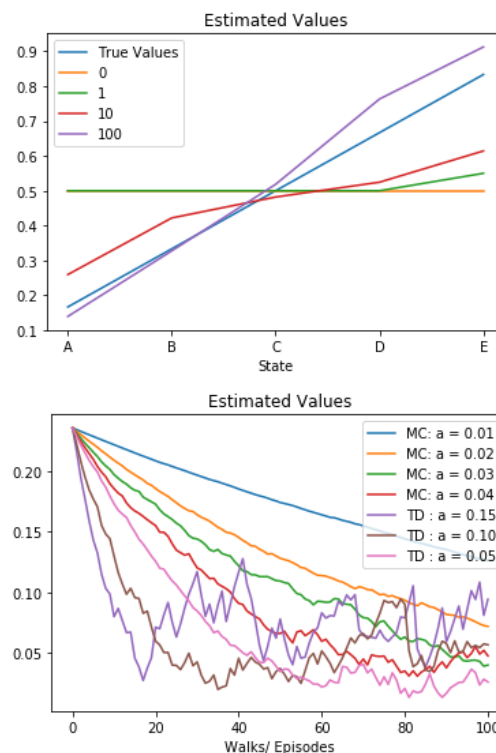
The problem of blackjack is modelled using two classes: Player and Dealer, and the problem is solved using First Visit MC. Since the player can HIT without losing whilst his sum is less than or equal to 11, the actual states which are of concern involve the sum from 12 to 21. Every ace is by default taken as 11 until the player/dealer go bust in which case the card is taken as 1. For better and diverse estimates, exploring starts has been used.

5. In the Driving home example, suppose we have accurate estimates for how long it takes to get Home ( $H_1$ ) from Work ( $W$ ), Now if one moves to a newer location ( $H_2$ ), we have the following:

From the past experience we know  $V(W)$ , now if we change the final destination to ( $H_2$ ), such that the initial journey still goes through the highway, then using Monte Carlo, we estimate  $V$  for all states  $s$ , by generating the entire episode again. Now since the updates are for the entire episode, the estimates for the initial part of the journey i.e. from  $W$  to Highway, also is recomputed, despite being the same.

When using Temporal Difference, the updates for each step, do not alter the estimates for the initial states (from  $W$  to Highway) which can essentially be used to estimate the value function  $V$ , for the newer states.

6. Attached in ipynb.



(6.3) The change happens only for one state, as the rewards for all other transitions is 0, with the reward being 0 for A, and 1 for E. and as the initial values are 0.5 the value can change only for one of A or E. also this change is precisely  $\alpha \cdot 0.5$ , that is 0.05. In this case, In the first episode the agent went to E, and thus got a reward of 1.

6.3 Starting at C, going to E

$$V(A) = 0.5$$

$$V(B) = 0.5$$

$$V(C) = 0.5 \quad (\text{initially})$$

$$V(C) = 0.5 + 0.1(0 + 0.5 - 0.5) = 0.5$$

$$V(D) = 0.5 + 0.1(0 + 0.5 - 0.5) = 0.5$$

$$V(E) = 0.5 + 0.1(1 + 0 - 0.5) = 0.5 + 0.05 = 0.55$$

thus  $V(E)$  changes by 0.55

(6.4) Trying out different values of alpha would definitely help in picking the better performing technique.

As the update:

$$V(S) = V(S) + \alpha(R + V(S') - V(S))$$

Is dependent on the step size, and for a smaller alpha the updates are less noisy and varied, while for an extremely smaller alpha the algorithm would take long to converge.

(6.5)

$$V(S) = V(S) + \alpha(R + V(S') - V(S))$$

As  $V(S)$  becomes to move closer to its actual value, the difference between the true value and the estimated value decreases, for a larger alpha, the step updates are larger which may cause  $V(S)$  to overshoot or diverge, thus increase the mean squared error.

7. Attached in ipynb.



8. Even when the action selection is made greedy, it is possible that the weight updates and the corresponding action sequences are different.

Let  $A_{\text{sarsa}}$  and  $A_Q$ , denote the actions picked by SARSA and Q-Learning respectively, for the state  $S_{t+1}$

In the case when the states  $S_t$  and  $S_{t+1}$  are different, then

$$A_{\text{SARSA}} = \operatorname{argmax}_a Q(S_{t+1}, a)$$

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R + Q(S_{t+1}, A_{\text{SARSA}}) - Q(S_t, A_t)]$$

$$A_Q = \operatorname{argmax}_a Q(S_{t+1}, a)$$

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R + \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

$$\text{Here } \max_a Q(S_{t+1}, a) = Q(S_{t+1}, A_{\text{SARSA}})$$

Thus the two updates and actions are the same

But when  $S_t$  and  $S_{t+1}$  are the same, then,

The action  $A_{\text{sarsa}}$  and  $A_Q$ , can be different, as  $A_{\text{sarsa}}$  is picked before the  $Q$  value is updated, i.e.  $Q(S_t, A_t)$ , while for Q-learning the action is picked in the beginning of the next step of the episode. This way  $A_Q$  is essentially the new greedy choice after the update, while  $A_{\text{SARSA}}$ , was the greedy choice before the update was made, and hence they may not be the same, and as a result the following state action sequence may be different leading to different weight updates as well.