# REPORT

**PHD17008**
**Aradhya Neeraj Mathur**

**Q1.**
**NOTE: the pseudocode has been used to implement the question 4 5.1 , 5.2.**
The pseudocode for question 4 is as follows
**G += R(t+1)**
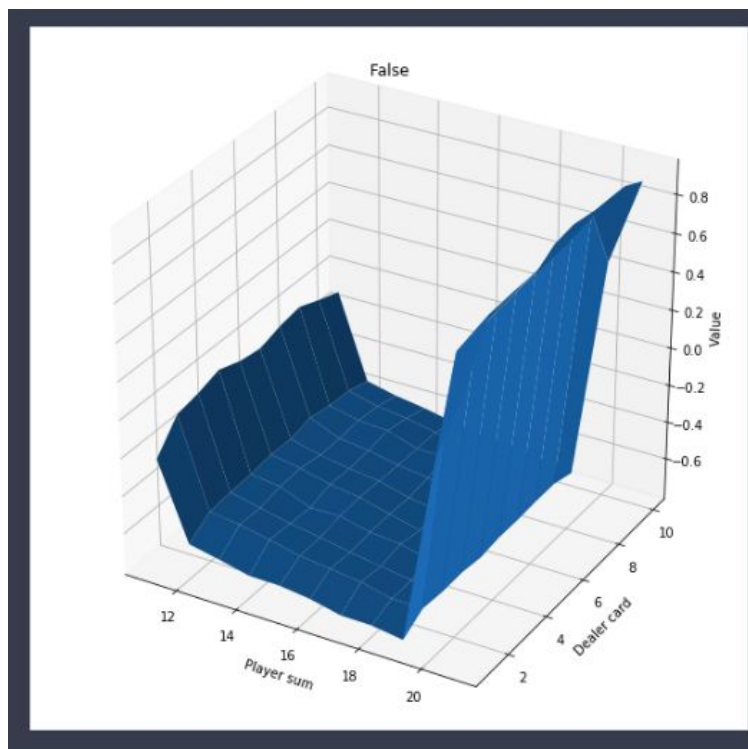**Q(s,  a) = (Q(s,  a) * count(S,a)  + G) / ( count(S,a) + 1)**
 While the pseudocode in book maintains a list of the returns and then calculates the average my method keeps track of the count of each state action pair and the multiplies the older Q value with the old count and adds the current updated return and then divides by the new count which is essentially equivalent to the mean.
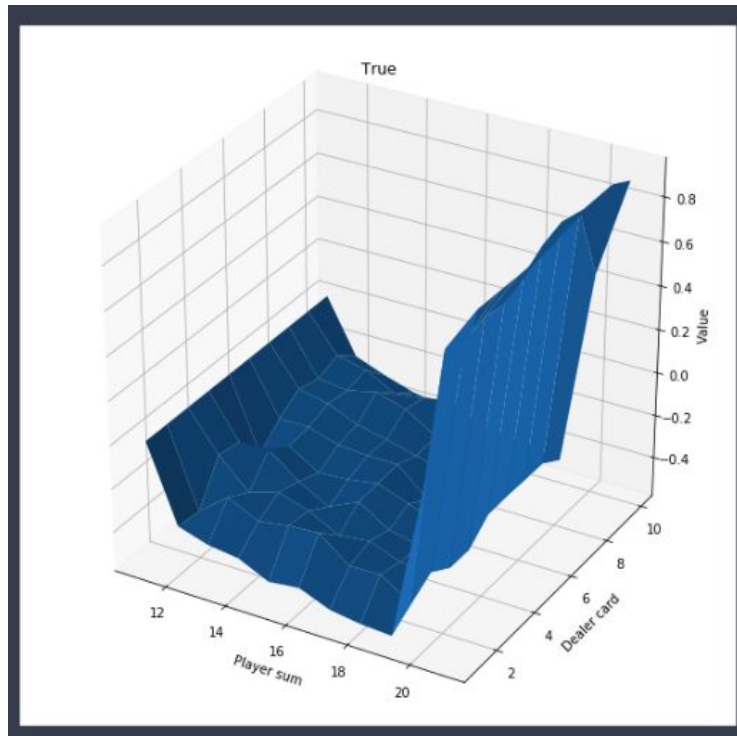
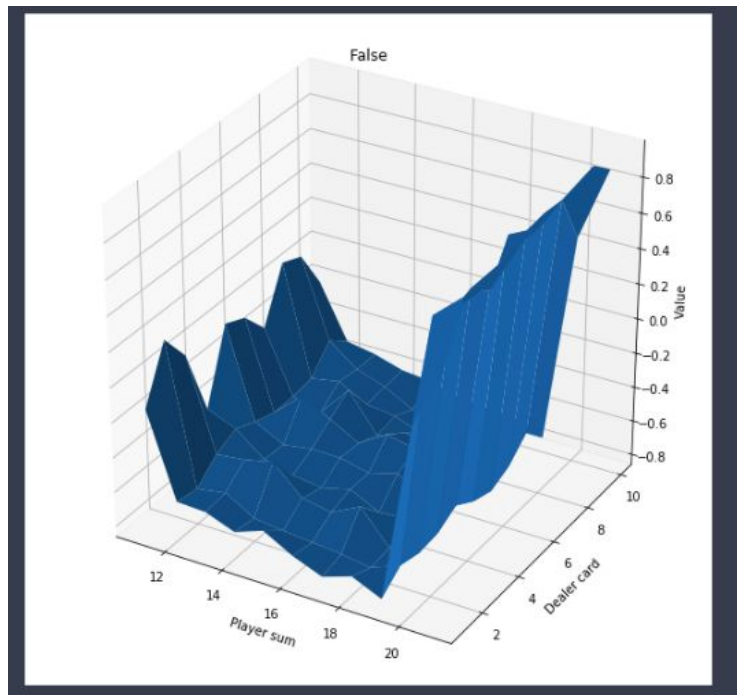**Q4 Value diagrams**
**Fig 5.1**
For 500000 episodes

**Non usable Ace**

**The code has been explained in the notebook**

## Usable Ace
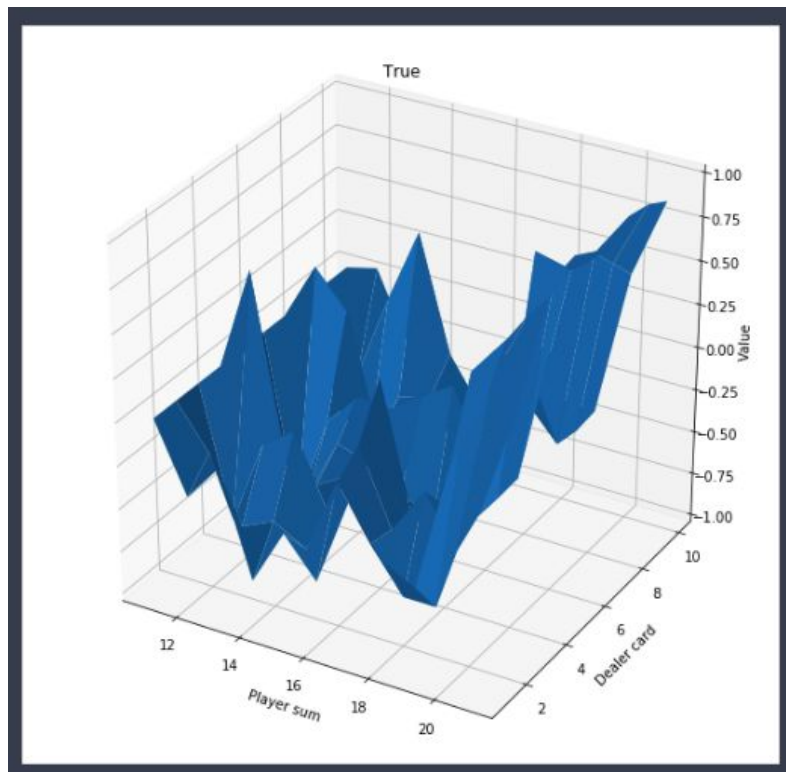


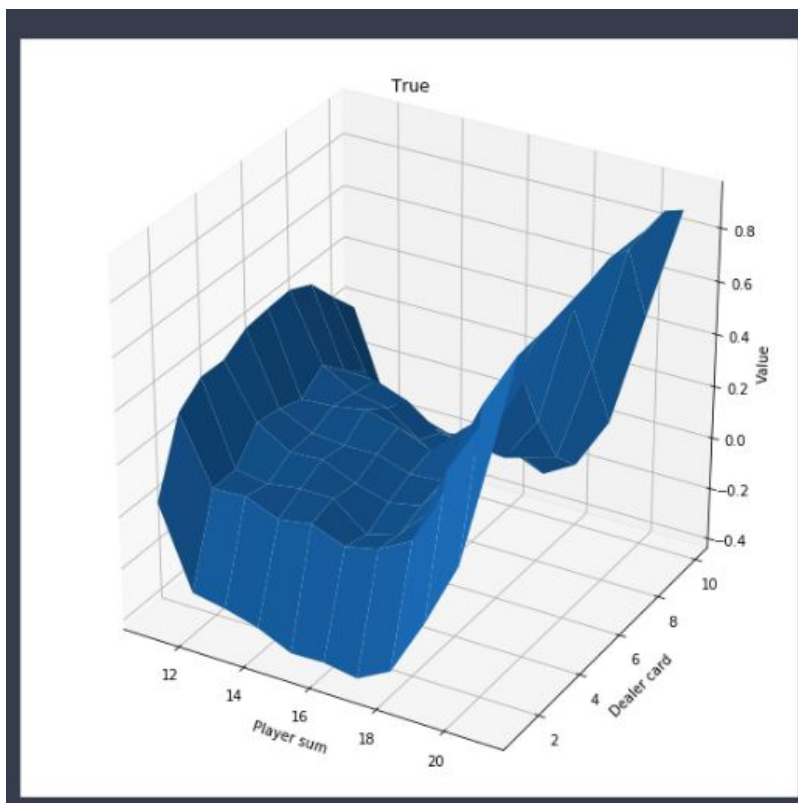**For 10000 episodes**

## Non usable Ace

**Usable Ace**
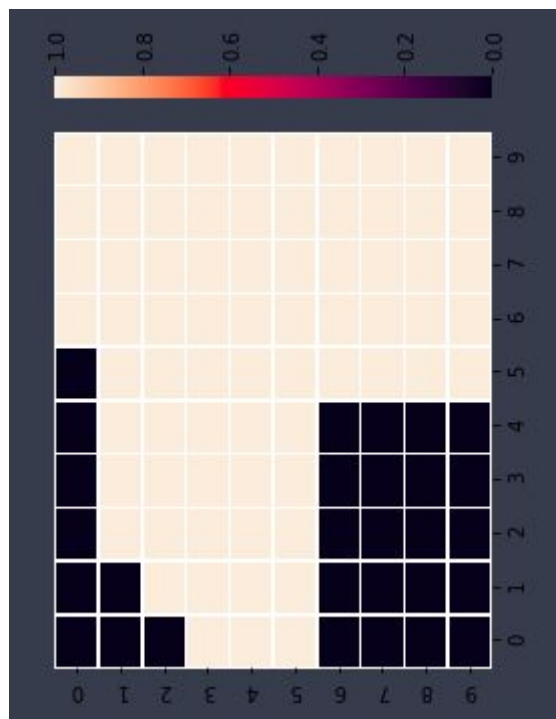


**Fig 5.2**

**Usable Ace**

## Non Usable Ace



## Policy diagram

## Non Usable Ace



**Non Usable Ace Policy Hit:0, Stick:1**

**Usable Ace Policy Hit:0, Stick:1**

**Fig 5.3**
**X axis :  number of episodes          y axis: MSE error in value**

**It can be observed from the graphs that the error of ordinary sampling is more than importance sampling.**

**Q6.**

**Empirical RMS error,**
**averaged over states for TD0 and MC for different alpha values**



**Q7.**



**Q Learning Vs sarsa for sum of rewards y axis:sum of rewards x axis:episodes**

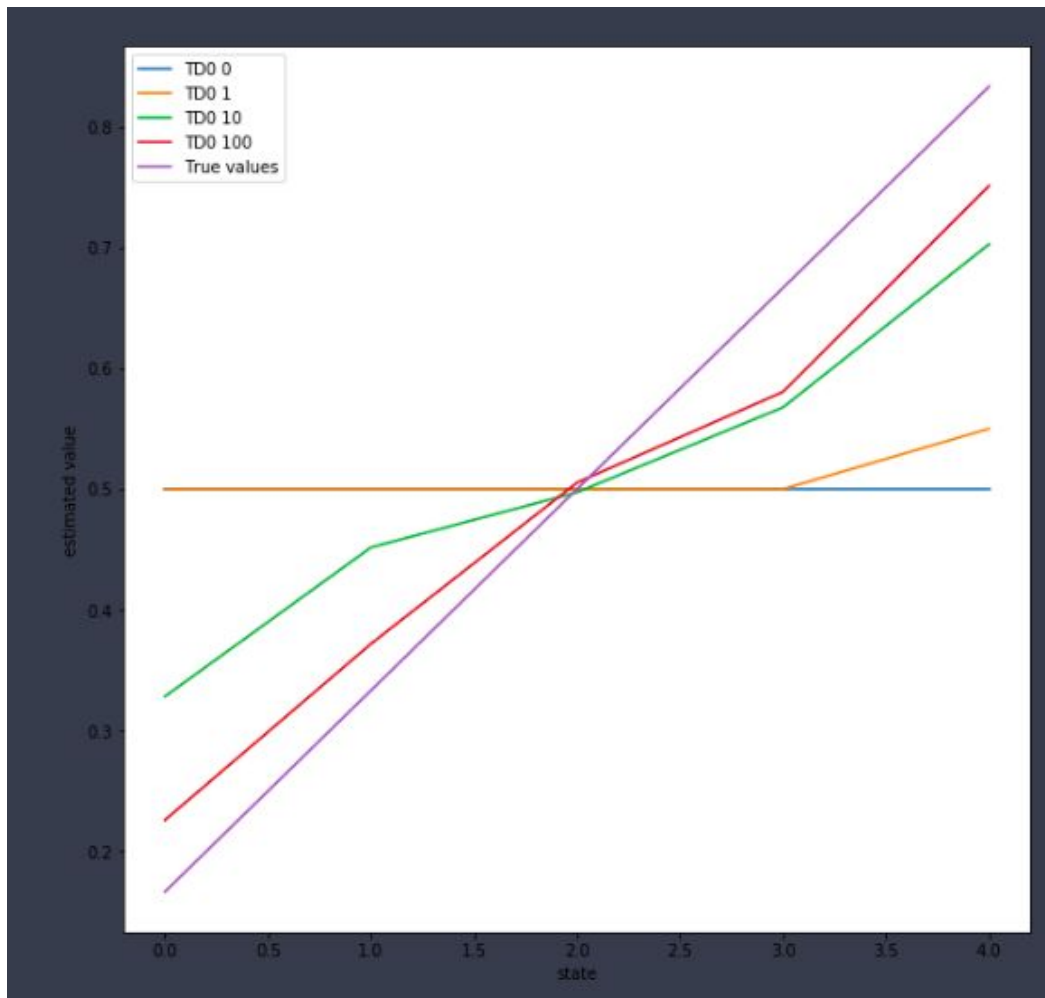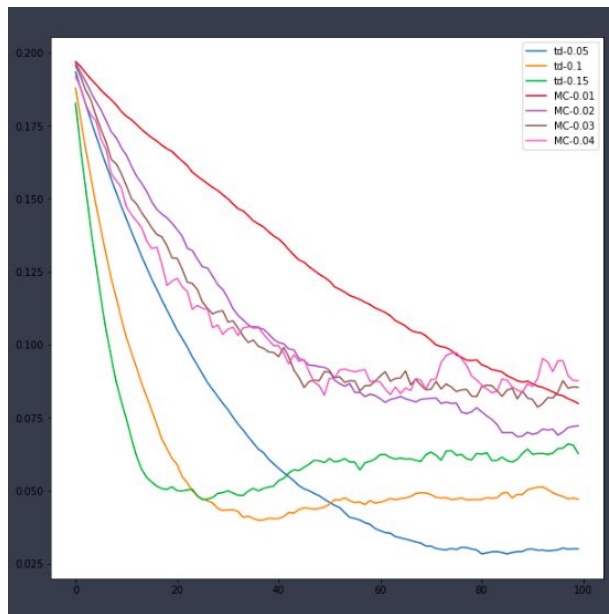equation 5.6 calculates weighted importance sampling using the following eqn.

$$V(s) = \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{T}(s)} \rho_{t:T(t)-1}}$$

Equation analogous for Q(s,a).

If we collect and cover the state action pairs in $\mathcal{T}(s,a)$ we can express Q(s,a)

as

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

$$\boxed{\rho_{t+1:T(t)-1} = \prod_{i=t+1}^{T-1} \frac{\pi(A_i|S_i)}{b(A_i|S_i)}}$$

T(t) is the first time of termination after time t and $G_t$ represents returns from time t+1 till t(t)

Ex. 5.5.

Backup diagram for 5.3 for $q_\pi$ is:

• $(s, a)$ — Starting with initialised 's'
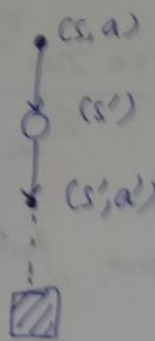from sequence of episodes
and action's corresponding
to state using policy $\pi$

○ $(s')$

× $(s', a')$

▨

Q1. Math for report explanation

$$Q_{n+1}(s,a) = \frac{1}{n}\left(\sum_{k=1}^{n} G_k\right) \Rightarrow \text{avg. of returns}$$

$$\therefore Q_{n+1}(s,a) = \frac{1}{n}\left[\sum_{k=1}^{n-1} G_k + G_n\right]$$

$$= \frac{1}{n}\left[\frac{(n-1)}{(n-1)}\sum_{k=1}^{n-1} G_k + \frac{1}{n}G_n\right]$$

$$= \frac{1}{n}\left[(n-1)Q_n + G_n\right] \rightarrow \text{Implemented in code}$$

$$Q_{n+1}(s,a) = Q_n(s,a) + \text{avg.}(\text{returns})$$

$$= Q_n(s,a) + \frac{1}{n}\left[\underset{\underset{\text{new count}}{\downarrow}}{(n-1)}\underset{\underset{\text{old count}}{\downarrow}}{Q_n} + \underset{\underset{\text{return}}{\overset{\rightarrow \text{current}}{}}}{G_n}\right]$$

$\therefore$ We keep track of the visits and current return
which is used to update the new $Q(s,a)$

**6.3:**

The TD method make the following update

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

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

we have $\gamma = 1$
$$\alpha = 0.1$$

In the first episode we only see a change the value of state A

Since $V(S) = 0.5$ for all $s \in S$

$$V(S') = 0 \text{ for } s' \text{ in terminal state}$$

now if in first episode we get reward 0 we end in the terminal state on the lyt.

$$V(A) = V(A) - 0.1 (V(A)) = 0.9 V(A)$$

$$\therefore V(A) = 0.9 \times 0.5 = 0.45$$

$$\text{difference} = 0.5 - 0.45$$
$$= 0.05$$

## 6.4 & 6.5

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

This $\alpha$ increases the weightage to the reward. The importance of each reward increases if a bigger value of $\alpha$ is used. Now as experimented in the notebook (plots can be found in the notebook) a smaller value of alpha leads to slower converge thus, if we keep the same number of episodes

we get poor performance. Though it might lead to smoother convergence. Any other alpha thus would not lead to better performance. If we increase alpha it will lead to more fluctuations. Similarly if α initialised to smaller values we

## Q8

### 5.12

For ~~TD~~ O-Q-learning

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

This modifies the state action value functions & then chooses the action as per update $Q(S, a)$.

~~Thoug~~ Whereas, in sarsa we choose action as per previous, state value function and then update it. Both lead to different evaluations.

## Q6.

### 6.2

In this case the part of bridge is common and only the initial route change. Assuming we use bootstrapping in TD, TD will be give a better performance as the state values of common states can be used which will mostly be common. This would give better convergence. Since there is less states that are arbitrary