

Ex-04

## Monte - Carlo Methods.

i) Incremental implementation of M-C methods.

First Visit MC prediction, for estimating  $V_{\pi}$  w.r.t. a policy  $\pi$  to be evaluated.

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$   
where  $N(s) \leftarrow 0$ .

$\therefore$  Loop Forever (for every episode)  
where  $G \leftarrow 0$

$\therefore$  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$   
 $G \leftarrow \gamma G + R_{t+1}$

Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$

$$N(S_t) \leftarrow N(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) + (G - V(S_t)) / N(S_t)$$

2) First visit vs every-visit.

a)

In the Blackjack task suppose every visit MC used instead of first-visit. The each episode is constantly changing which will appear only once. Even if every visit MC is used. Since the state appears only once, it gets the same result as using first-visit.

b) MDP with a single nonterminal state and a single action that transitions back to the non-terminal state with probability  $p$  and terminal state  $1-p$ .

$$\begin{matrix} T=10 \\ \gamma=1 \end{matrix} \quad \left\{ \begin{matrix} \text{(given)} \end{matrix} \right.$$

$$G_t = R_t + G_1 = 10$$

For First-visit

$$G_1 = R_2 + G_2 = 9$$

$$V(s) = 10$$

$$G_2 = R_3 + G_3 = 8$$

$\vdots$

For every visit

$$G_9 = R_{10} + G_{10} = 1$$

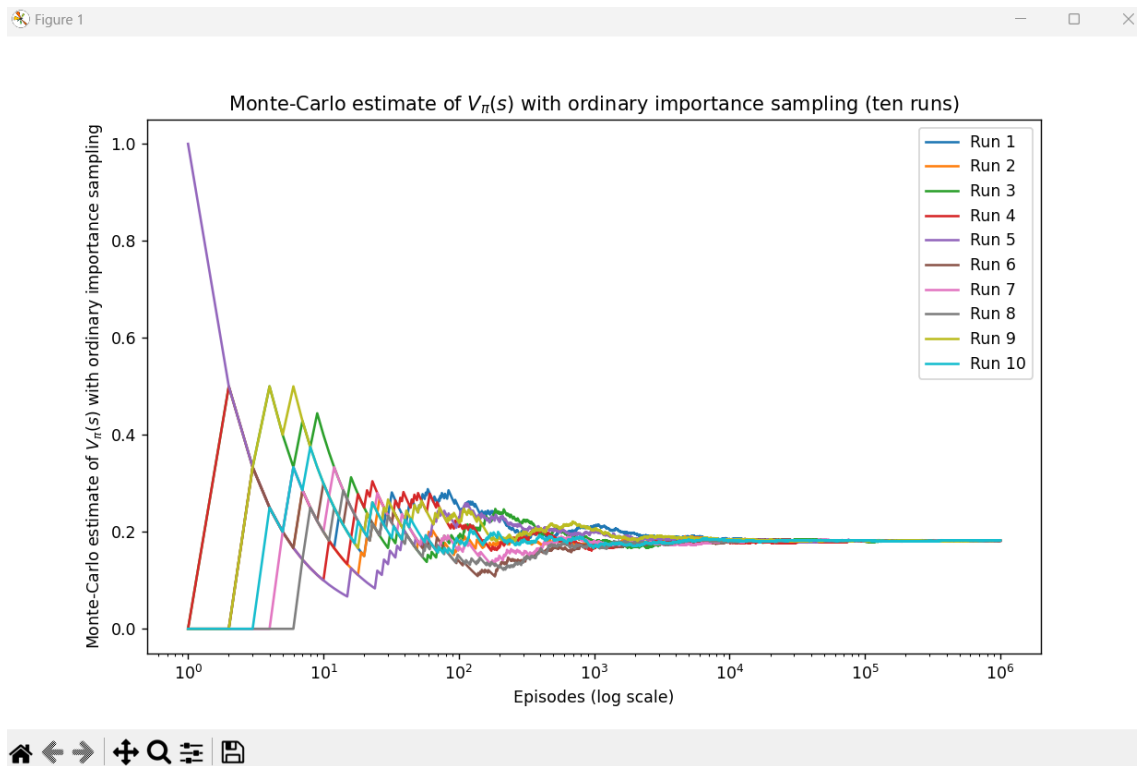
$$V(s) = \frac{G_{10} + G_9 + \dots + G_0}{10}$$

$$G_{10} =$$

$$= \frac{0 + 1 + \dots + 10}{10}$$

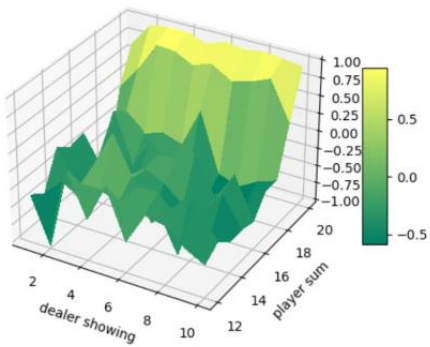
$$V(s) = 5.5$$

2(c)

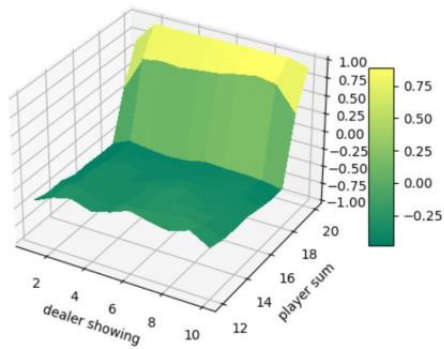


3) a), b)

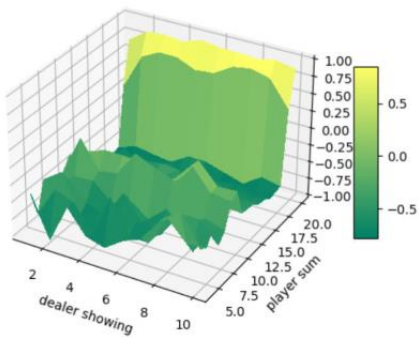
Usable Ace after 10000 episodes



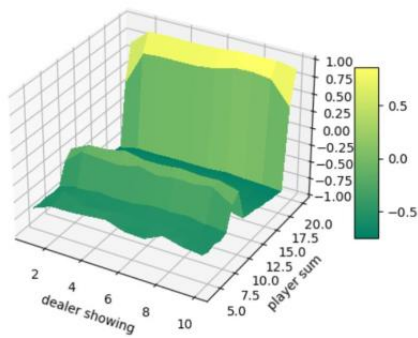
Usable Ace after 500000 episodes



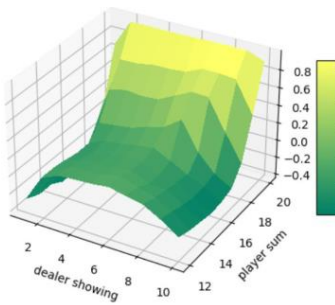
No usable Ace after 10000 episodes



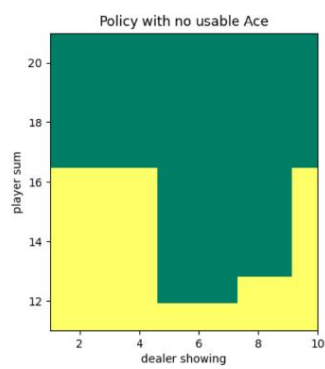
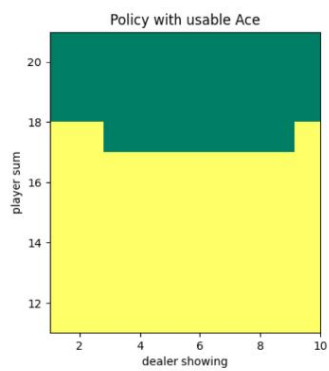
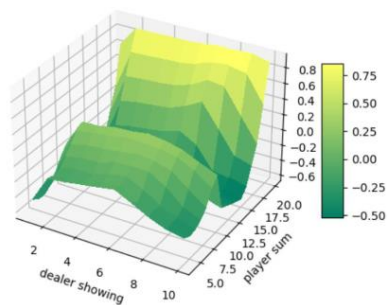
No usable Ace after 500000 episodes



Usable Ace after 5000000 episodes

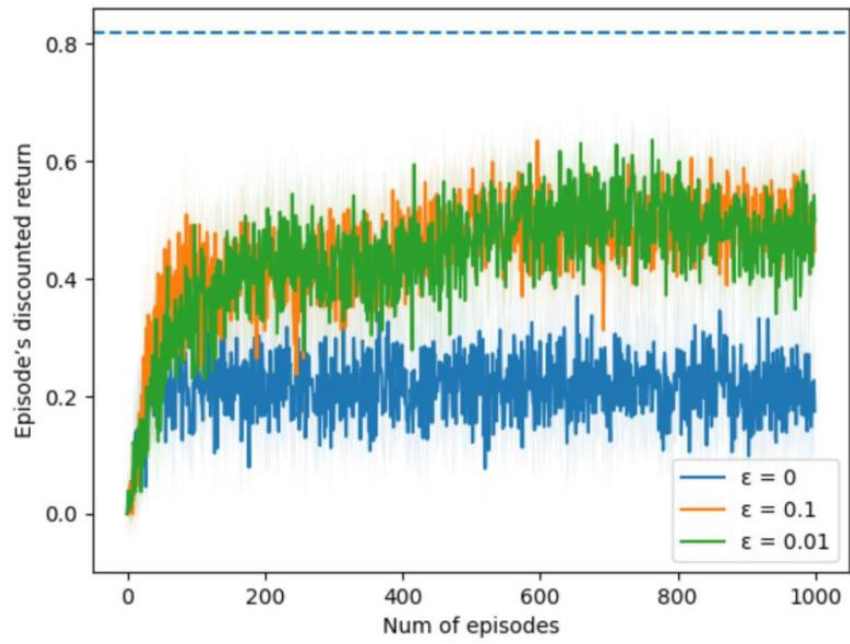


No usable Ace after 5000000 episodes



4)

b)





4) Four Rooms, revisited.

c) Importance of doing exploring starts in M-CES.  
when  $G=0$  ~~get~~ without exploring start point will always choose the state with the 1<sup>st</sup> the result & refuse to make a new attempt. However with Exploring start this is prevented. It will have chances to try other states. Thus this will still be useful for improving the policy.