

RL HW EX4:

Question 1:

- a. The results would be pretty much the same. This is because in this blackjack example there is only one way that the same state can be reached twice in the same episode is if an ace is flipped from 11 to 1 and you end up with a total value that is the same as you have already encountered. This is a rare case and thus wouldn't have much of an effect on the state-value functions.

- b. First visit:

$$G = 10/1 = 10$$

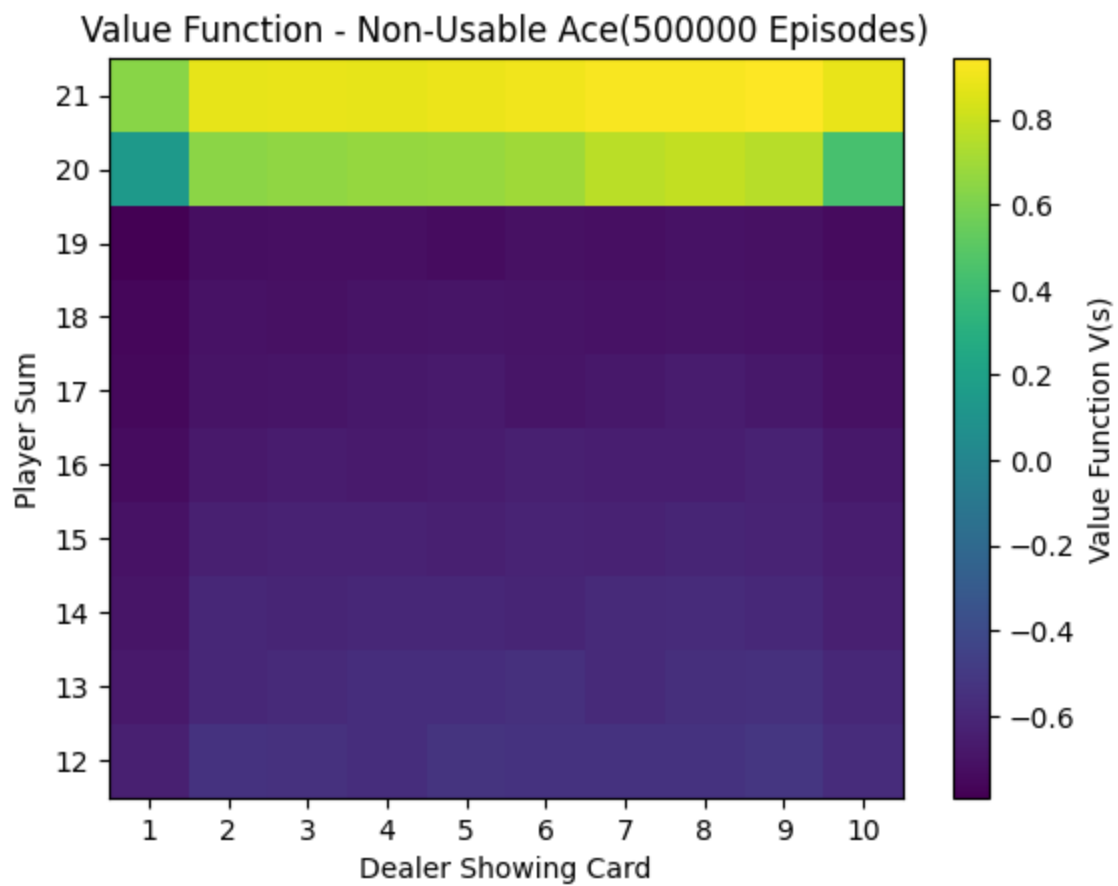
Every visit:

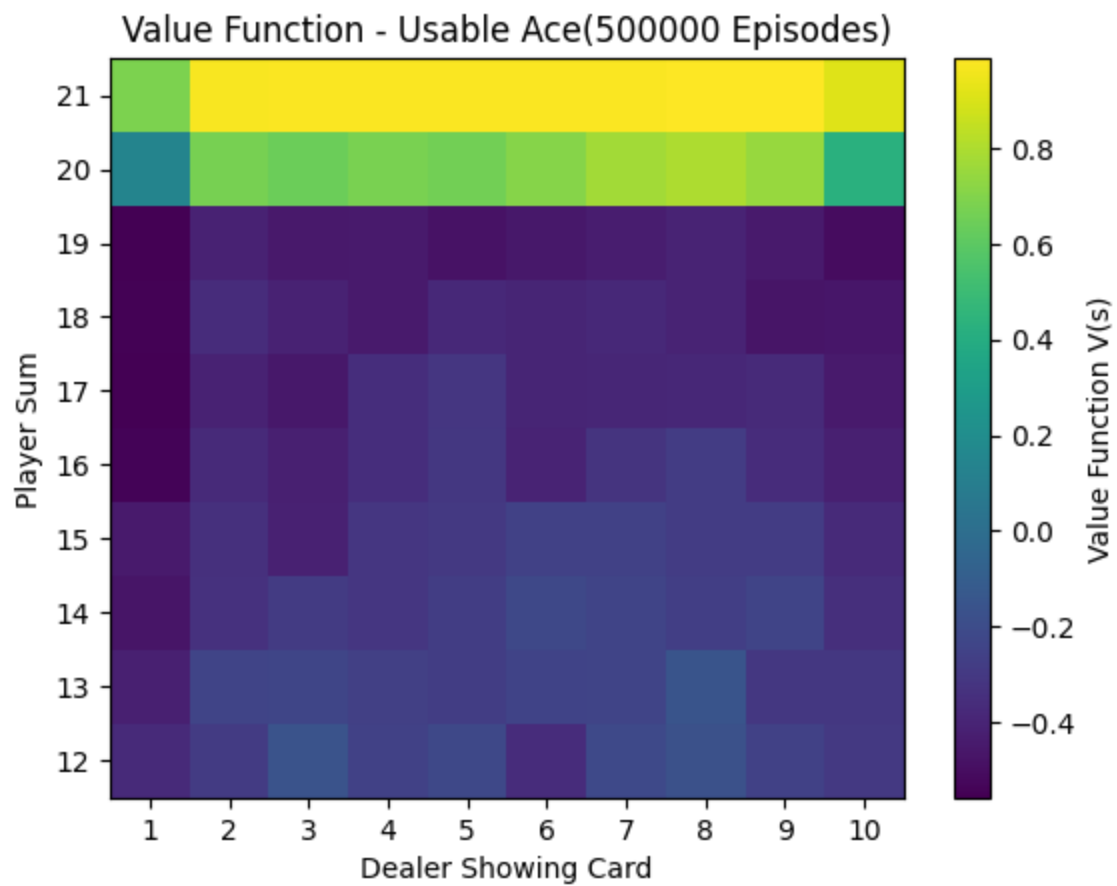
$$G = \text{avg}(1,2,3,4,5,6,7,8,9,10)$$

$$G = 5.5$$

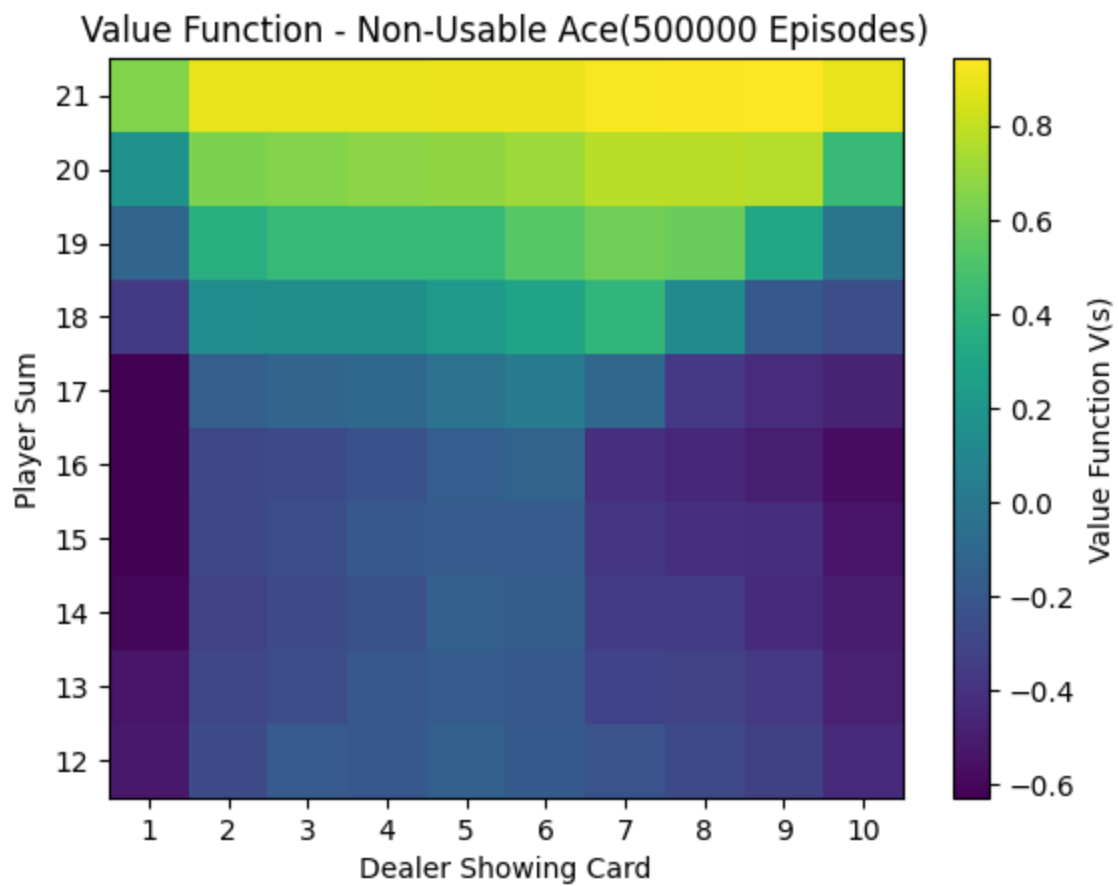
Question 2:

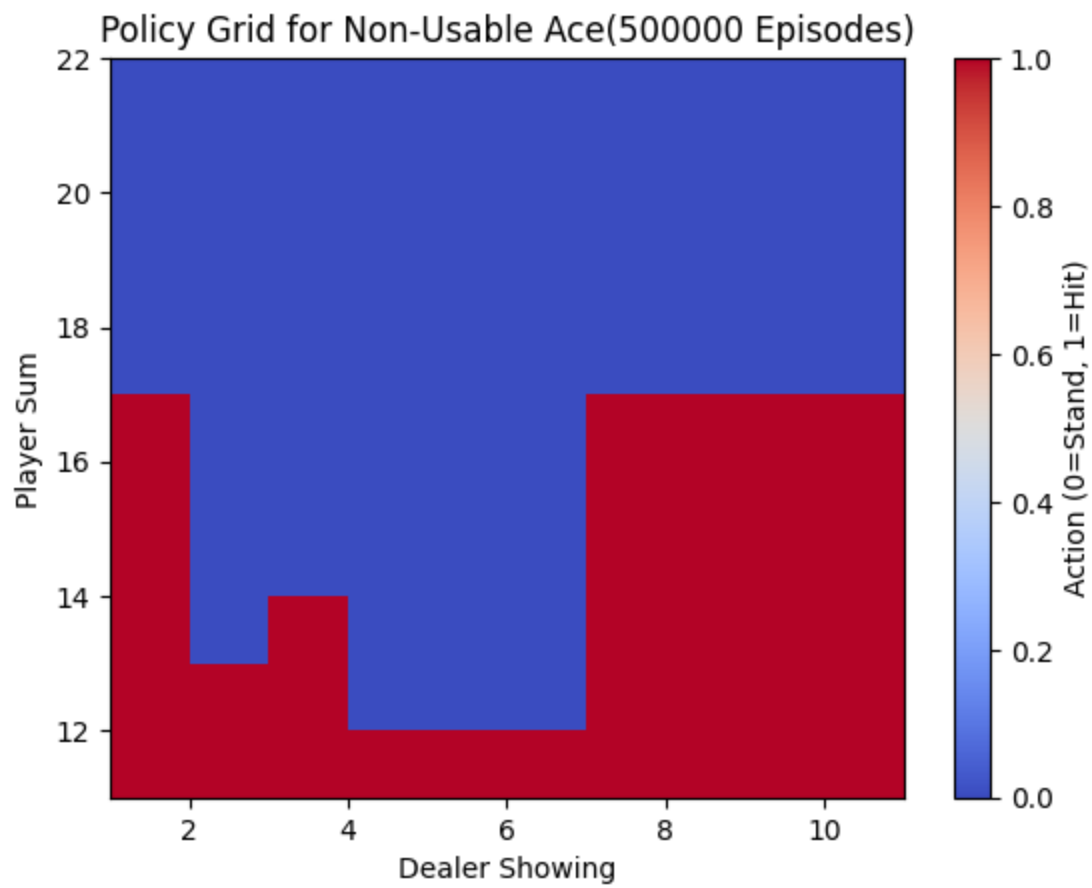
- a.

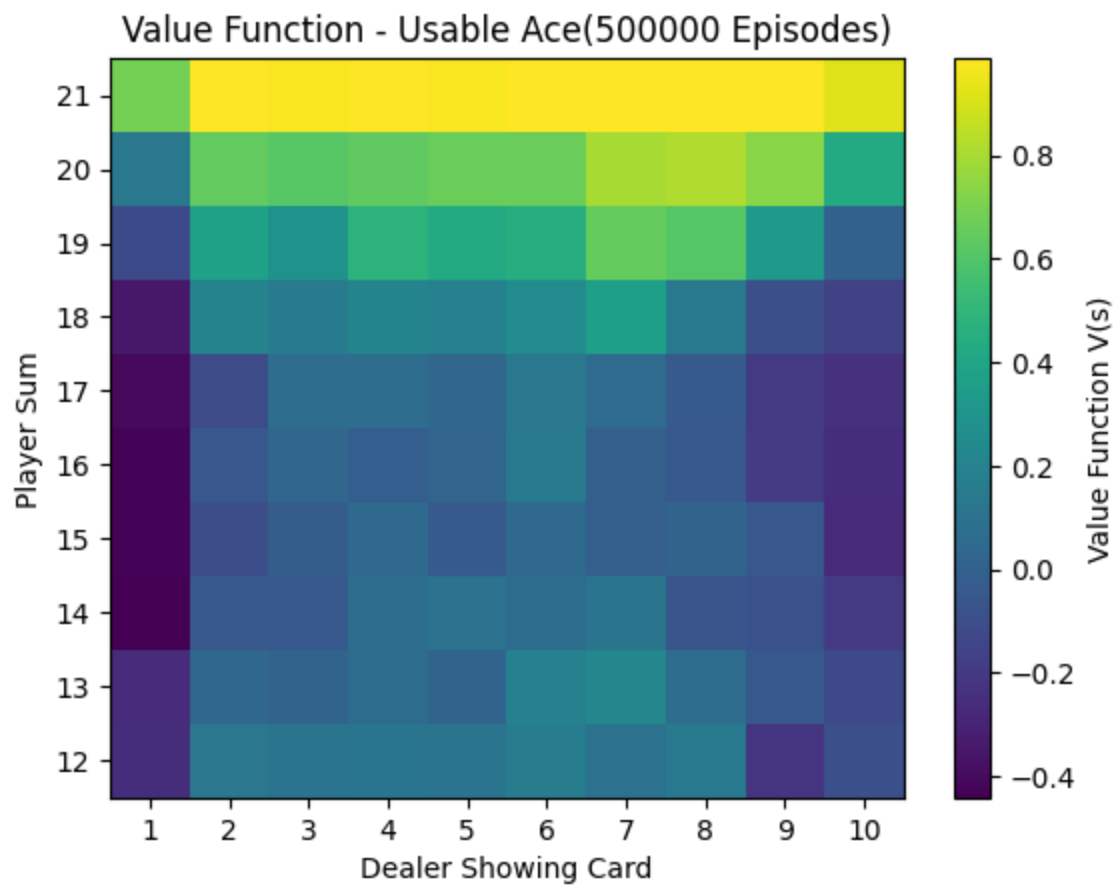


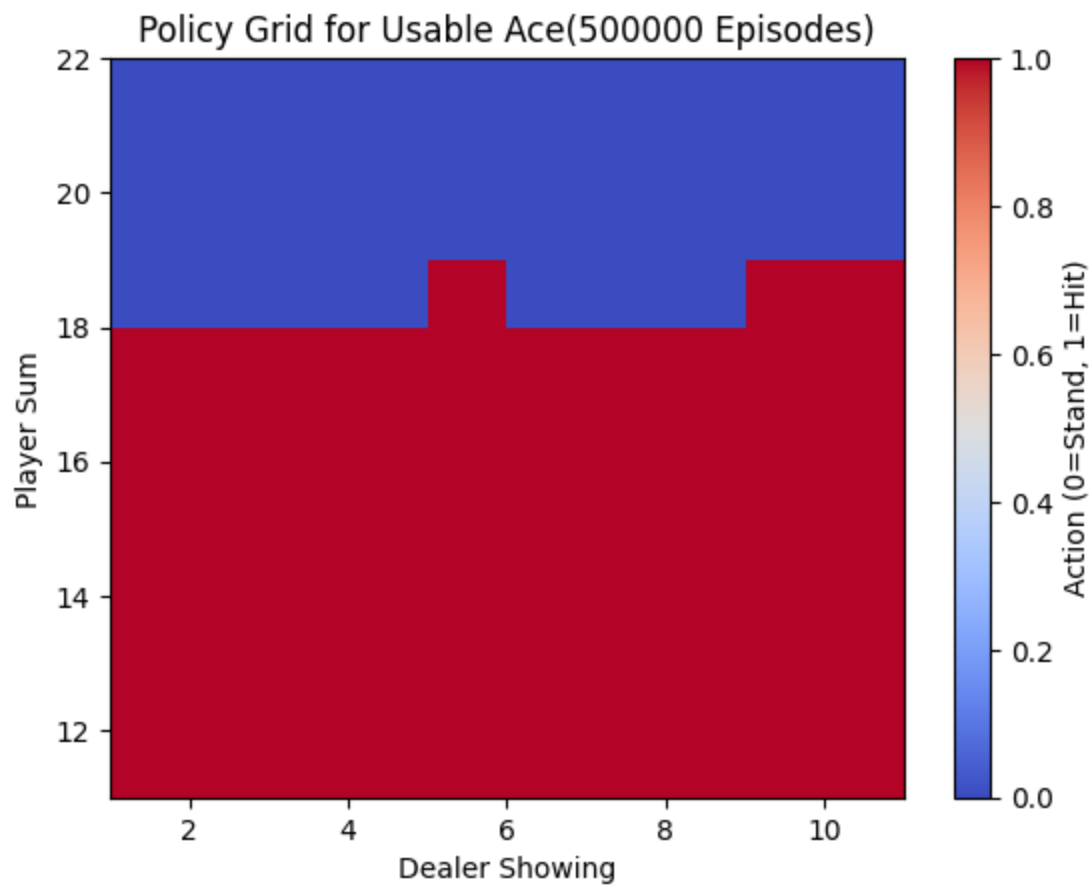


b.



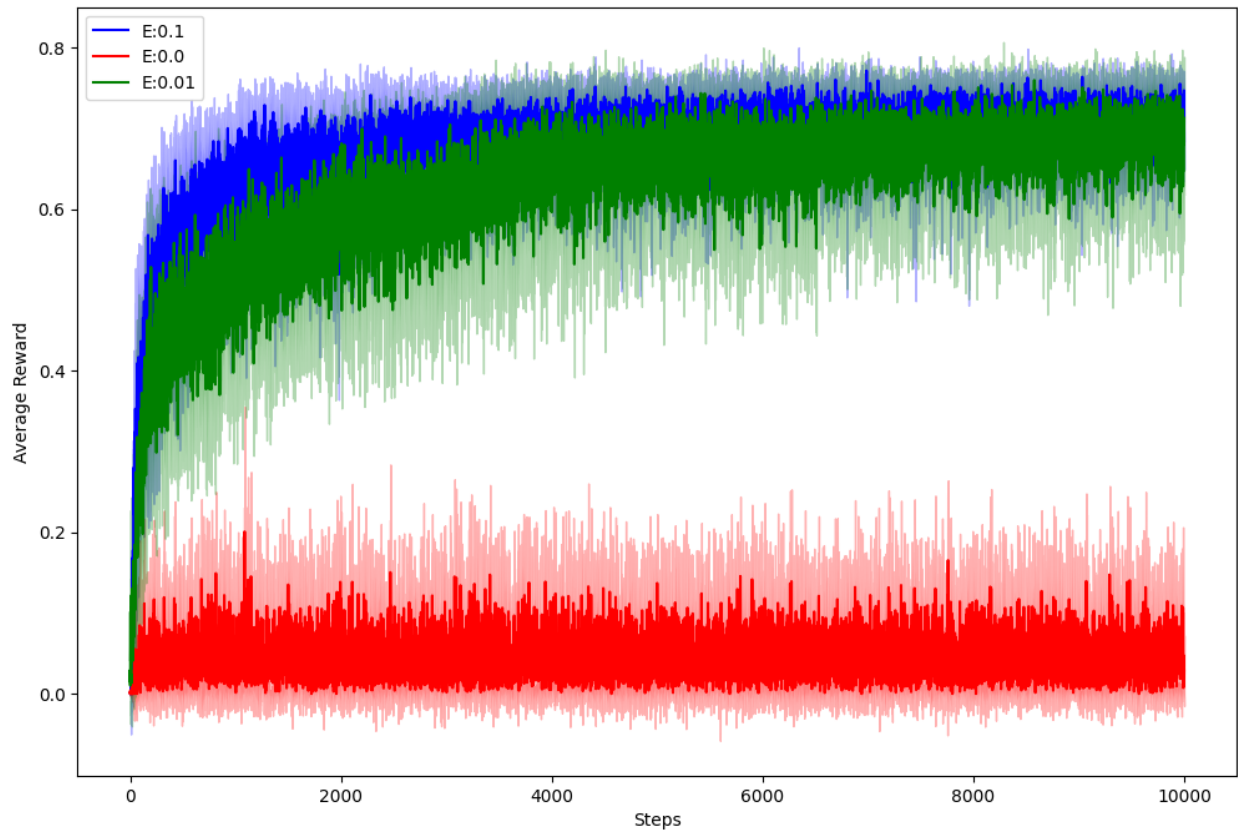






Question 3:

a.



- a. The results for $\epsilon = 0$, highlight the importance of doing exploring starts. It highlights this as exploring starts ensure optimality by guaranteeing that all state-action pairs are explored, which lets the agent to accurately assess the value of each action in every state. This exploration is necessary because, without trying all possible actions, the agent cannot know if its current policy is truly optimal. Similarly, an ϵ -greedy policy ensures exploration by occasionally selecting random actions, helping avoid getting stuck in suboptimal behavior. If neither exploring starts nor ϵ -greedy exploration are used, the agent will only exploit its current knowledge and may never discover better actions, meaning it can never guarantee reaching the optimal policy.

Question 4:

- b. The policy π is a greedy and deterministic policy, we only observe the trajectories where the probability of action A_t is selected in state S_t is 1(

$\pi(A_t|S_t) = 1$) so $1/b(A_t, S_t)$ is correct for estimating the importance-sampling ratio in this case.