

Please refer readme.txt before running the code.

Answer for Q1:

1. Initialize

$$V(s) \leftarrow 0, \text{ for all } s \in S$$

$$N(s) \leftarrow 0, \text{ for all } s \in S$$

Loop forever :

Generate episode following π : $s_0, a_0, R_1, s_1, a_1, \dots, R_T$

$$G \leftarrow 0$$

Loop $t = T-1, T-2, \dots, 0$:

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

Unless s_t appears in s_0, s_1, \dots, s_{t+1} :

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

$$V(s_t) \leftarrow V(s_t) + 1/N(s_t) (G - V(s_t))$$

1. b. Initialize: $Q(s, a) = 0$

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

$$\pi(s) \in A(s) \text{ arbitrarily}$$

Loop forever :

Choose $s_0 \in S, a_0 \in A(s_0)$ randomly

Generate an episode from s_0, a_0 , ~~following~~ following π : $s_0, a_0, R_1, \dots, s_{T-1}, a_{T-1}, R_T$.

$$G \leftarrow 0$$

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

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

Unless the pair s_t, a_t appears in s_0, a_0, \dots, a_{t+1} :

$$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + 1/N(s_t, a_t) (G - Q(s_t, a_t))$$

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

Answer for Q2:

2. (a) No different: ~~Each ep~~
Each state only appears one time in the episode.

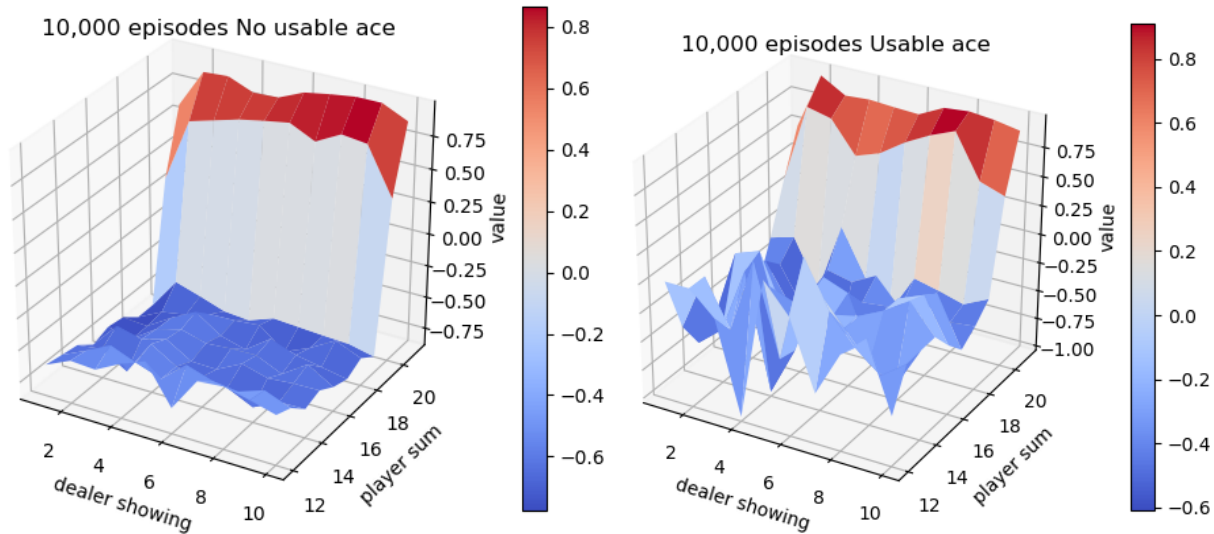
(b) $S_1 A_1 R_1, S_2 A_2 R_2, S_3 A_3 R_3, \dots, S_{10} A_{10} R_{10}$

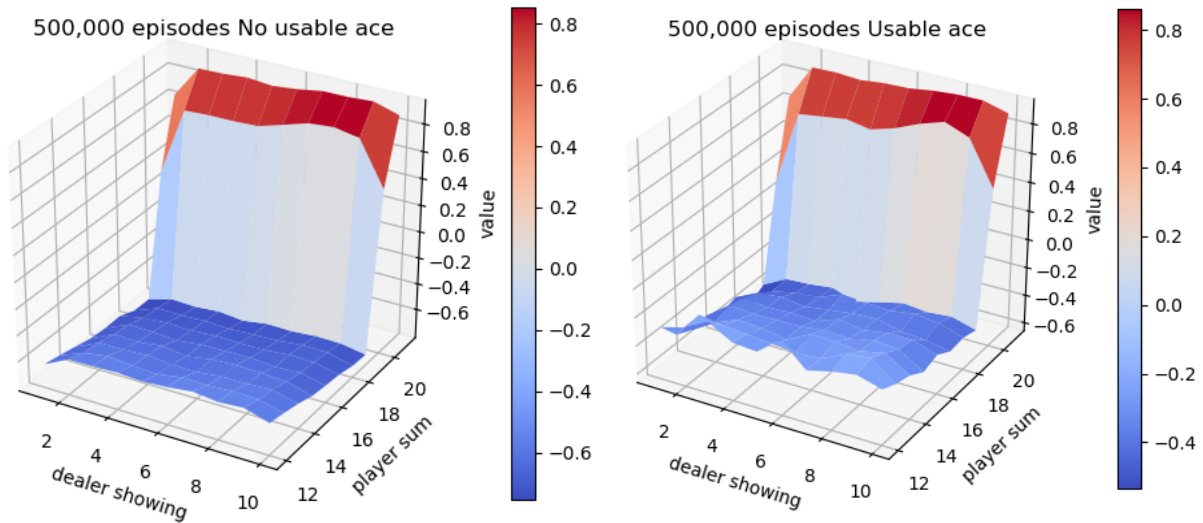
First visit: $V(s) = G_{1:T} = 10$

Every visit: $V(s) = \frac{G_{1:T} + G_{2:T} + \dots + G_{T:T}}{10} = \frac{55}{10} = 5.5$

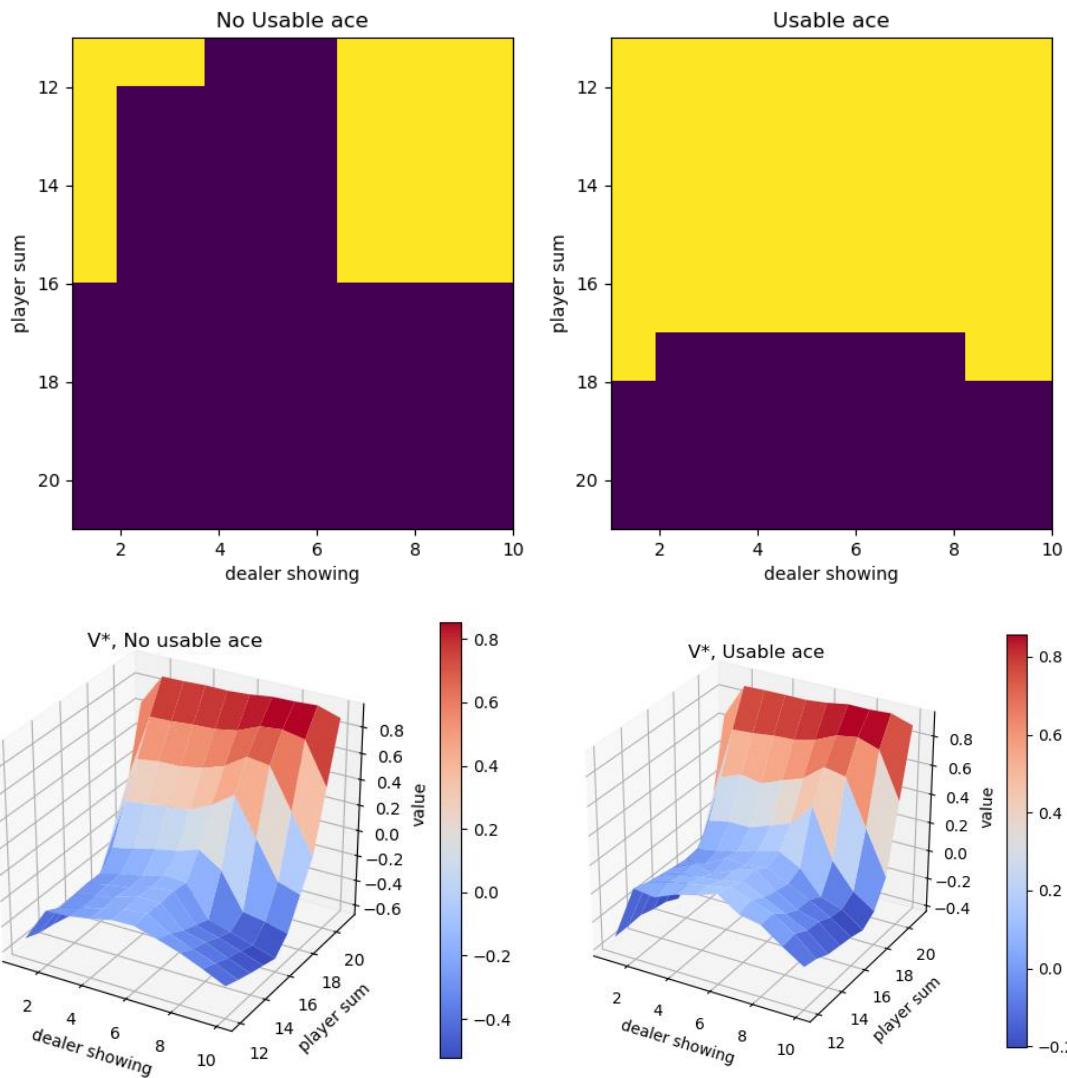
Answer for Q3:

- a) run python q3_a.py to generate the data (the default episode number is 10,000)
run python q3_a_plot.py to plot the figures below.



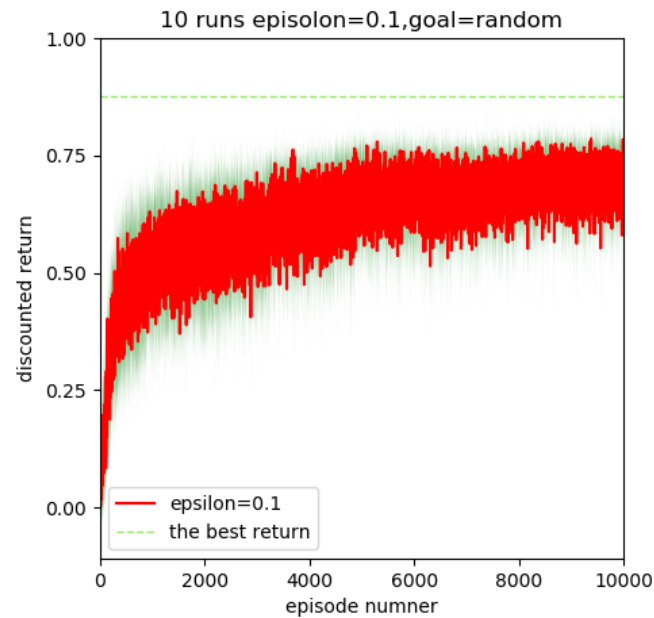
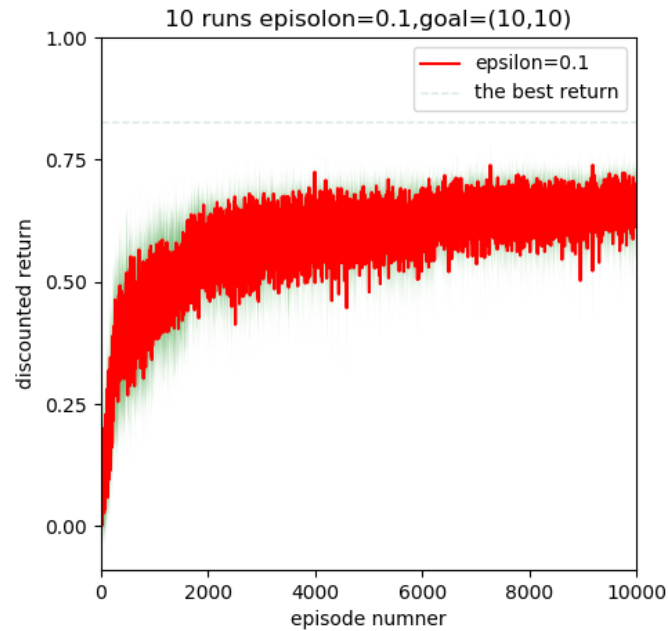


b) python q3_b.py to generate the data; python q3_b_plot is used to plot the figure below



Answer for Q4:

- a) python q4_a.py is used to generate the data, python q4_b is used to plot the figure below:

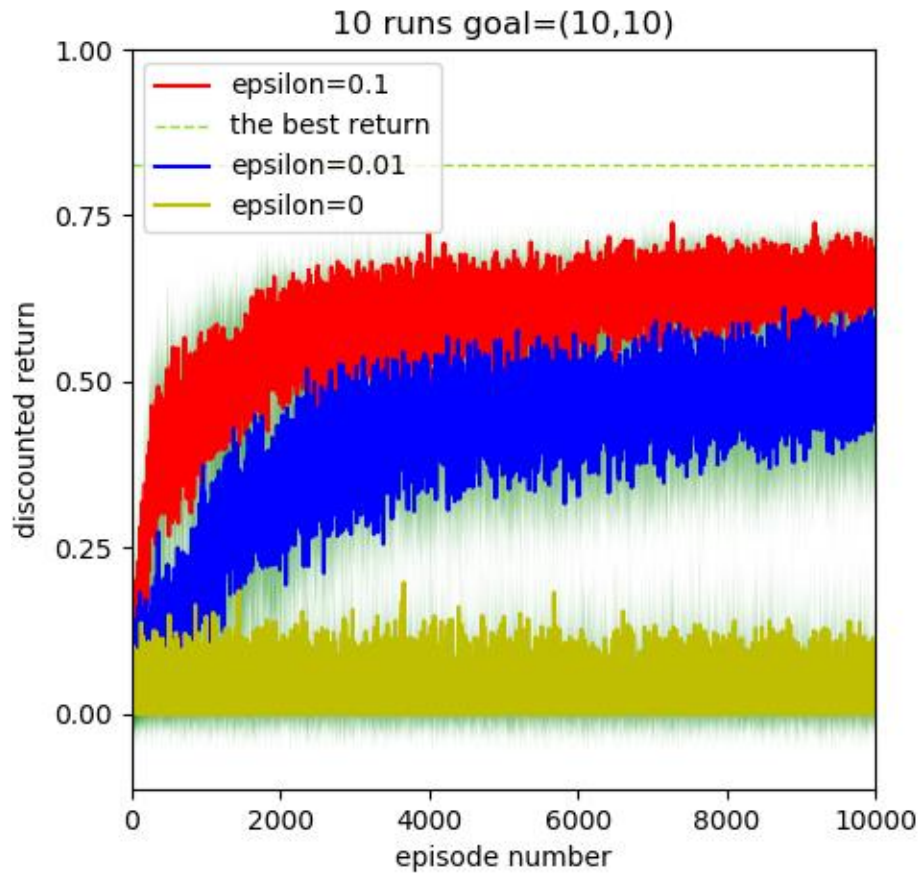


Details: the figures above are ten runs with episode number of 10,000, one for goal position (10,10), one for random goal position. The initial policy is random. The green fade area are the confidence bands.

- With first-visit MC, the agent learns to reach the goal.

b)

q4_a.py with the specific parameters is used to generate the data; q4_b_plot.py is used to plot the figure below:



The confidence bands are in shade green colors

c)

when $\epsilon=0$, there is no explore and the agent learns nothing. But in MCES, the random start state and start action is used to explore and substitute the epsilon greedy method.

Answer for Q5:

5. (a) given:

$$V_n = \frac{\sum_{k=1}^{n-1} W_k \cdot G_k}{\sum_{k=1}^{n-1} W_k}, \quad n \geq 2$$

$$C_{n+1} = C_n + W_{n+1}, \quad C_0 \doteq 0$$

prove

$$\therefore V_n = \frac{\sum_{k=1}^{n-1} W_k \cdot G_k}{C_{n-1}} \rightarrow V_{n+1} = \frac{\sum_{k=1}^n W_k \cdot G_k}{C_n}$$

$$V_{n+1} = \frac{(C_{n-1} + W_n - W_n) \cdot \left(\sum_{k=1}^{n-1} W_k \cdot G_k + W_n \cdot G_n \right)}{C_n},$$

$$= \frac{(\cancel{C_{n-1}} + W_n) \cdot \sum_{k=1}^{n-1} W_k \cdot G_k}{C_n \cdot C_{n-1}} + \frac{-W_n \cdot \sum_{k=1}^{n-1} W_k \cdot G_k}{C_n \cdot C_{n-1}} + \frac{W_n \cdot G_n}{C_n}$$

$$= \frac{\sum_{k=1}^{n-1} W_k \cdot G_k}{C_{n-1}} + \frac{W_n}{C_n} \left(C_n - \frac{\sum_{k=1}^{n-1} W_k \cdot G_k}{C_{n-1}} \right)$$

$$= V_n + \frac{W_n}{C_n} (C_n - V_n), \quad n \geq 1$$

When ~~$n \geq 1$~~ , $n=1$,

$$V_2 = V_1 + \frac{w_1}{c_1} [q - V_1]$$

$$\therefore C_1 = C_0 + w_1 = w_1$$

$$\therefore V_2 = V_1 + q - V_1 = q$$

$\therefore V_2$ doesn't depend on V_1 , so V_1 could be arbitrary

$$\begin{cases} V_{n+1} = V_n + \frac{w_n}{c_n} [C_n - V_n] \\ C_{n+1} = C_n + w_{n+1} \end{cases}, \left(n \geq 1, C_0 = 0, \right. \\ \left. V_1 \text{ is arbitrary} \right)$$

b). Target policy is greedy. It assigns the action
with largest q value the probability 1.

Answer for Q6:

a)

q6_a.py is used to generate the data and saved the trajectory with probability of action taken in the file named 'episode.p'. It also saves the final policy and the q values.

b)

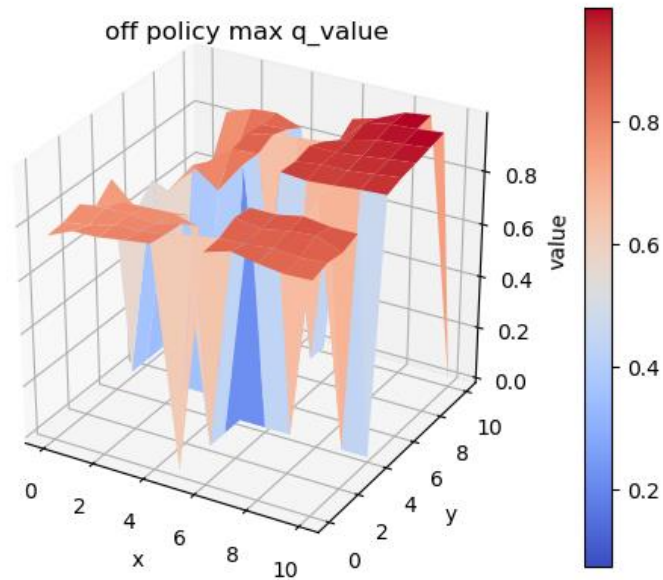
q6_b.py is used to compute the greedy policy respect to the q value stored in part (a). The greedy policy is shown below:

down	right	left	left	right	wall	right	up	right	right	right
right	down	up	down	right	wall	left	down	right	right	up
left	right	up	right	right	right	down	right	right	right	up
left	right	down	up	right	wall	right	right	right	up	up
right	right	down	right	right	wall	up	up	right	up	up
wall	left	wall	wall	wall	wall	right	right	up	up	left
right	left	up	down	up	wall	wall	wall	up	wall	wall
right	right	down	right	down	wall	right	right	up	left	left
up	up	right	down	right	wall	right	right	up	up	up
right	right	down	right	right	right	right	up	up	up	up
up	left	right	right	up	wall	up	right	up	up	up

- The greedy policy is better than the policy generated in the part (a) but it has no exploration and it doesn't converge to the optimal policy.
- The policy is suboptimal or gives some bad actions in some cells like (4,10), (3,10)

c)

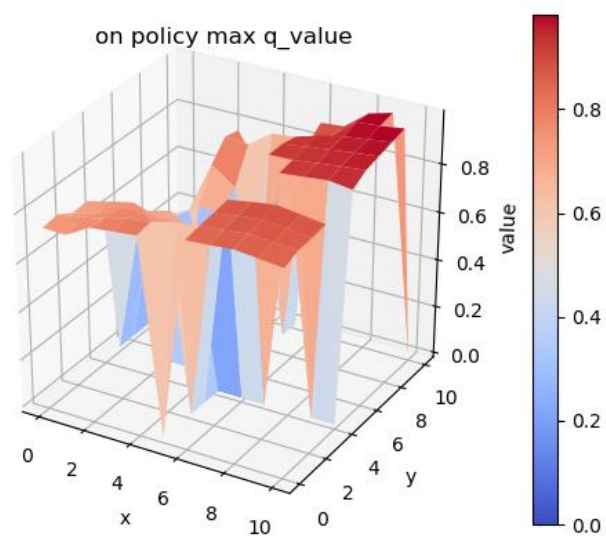
q6_c.py is used to do the off-policy prediction, instead of plot the q values for each action pair. I plot the max q value for each state as below:



Where the initial position is in (0,0) and the goal position is in (10,10). The value of the wall is set to zero.

d)

the code q6_d.py is used to calculate the q value for the greedy policy and q6_d_plot.py is used to plot the max q value for each cell.



e)

- Off-policy prediction value seems coarser (some area lower, some area higher) than the on-policy prediction due to the variance
- Off-policy method seems to converge to the on-policy prediction
- Due to the explore of the behavior policy, the target policy could learn the value of uncommon states. (the top left room)