1.a)

Initialize:

$\qquad V(S) \in \mathbb{R}$, arbitrarily $\forall S \in \mathbb{N}^+$

$\qquad C(S) = 0 \quad S \in \mathbb{N}^+$

Loop for ever ( each episode)

$\qquad$ Generate an episode $\pi$: $S_0 - - - R_T$

$\qquad G \leftarrow 0$

$\qquad$ loop for each step $t = T-1, T-2, --- 0$

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

$\qquad\qquad$ unless $S_t$ appears in $S_0 - S_{t+1}$

$\qquad\qquad\qquad C(S_t) \leftarrow C(S_t) + 1$

$\qquad\qquad\qquad V(S_t) \leftarrow V(S_t) + \dfrac{1}{C(S_t)} (G - V(S_t))$

b) Replace Initialization of return map with

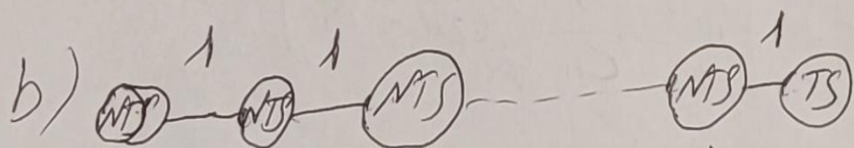$\qquad C(S,a) \leftarrow 0 \quad \forall S \in \mathbb{N}^+, a \in A(S)$

Replace "Append G in Return" with

$\qquad C(S_t, a_t) \leftarrow C(S_t, a_t) + 1$

Replace "$Q \leftarrow average( \;,$" with

$\qquad Q(S_t, a_t) \leftarrow Q(S_t, a_t) + \dfrac{1}{C(S_t, a_t)} (G - Q(S_t, a_t))$

2-a) No, by the nature of blackjack we will essentially see a state once in each episode so first-visit and every-visit make the same result.

b)



$\Rightarrow$ first-visit: $U(s, NTS) = \sum_{i=1}^{10} 1 = 10$

every-visit: $V(s, NTS) = \frac{1}{10}(1+2+3+\cdots 10) = \frac{1}{10}\sum_{i=1}^{10} i = 5.5$

c) Yes, according to the book which says in

$V[X] = E[X^2] - \bar{X}^2$  if $E[X^2] \rightarrow \infty \Rightarrow V[X] \rightarrow \infty$

$\bar{X}^2$: finit

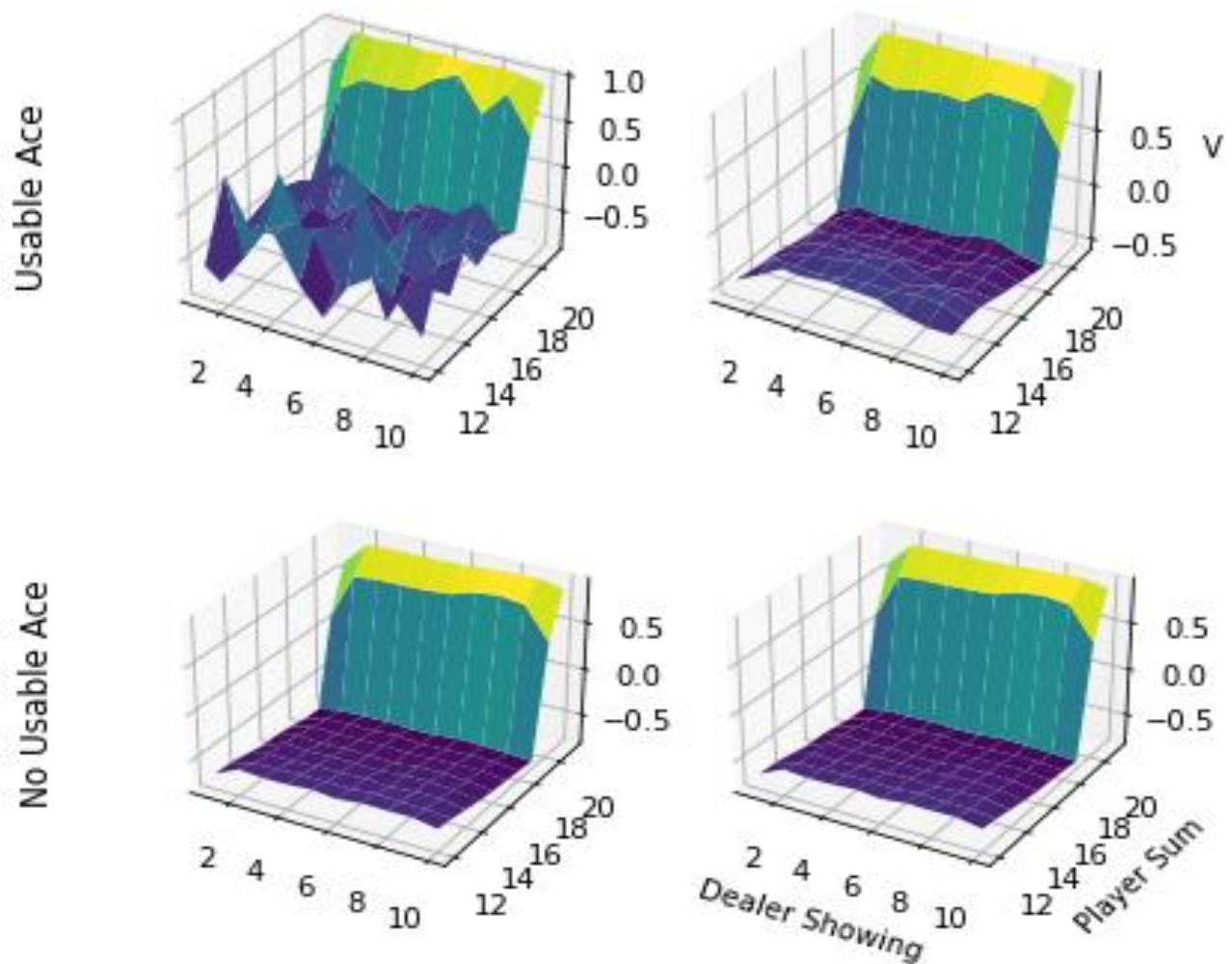for first-visit: $E_b\left[\left(\prod_{t=0}^{T-1} \frac{\pi}{b} G_0\right)^2\right] \rightarrow \infty$   we know

$V(X) \rightarrow \infty$ for every-visit

for every-visit we will have:

$E_b\left[\left(\frac{1}{T-1}\sum_{K=1}^{T-1} \prod_{t=0}^{K} \frac{\pi}{b} G_0\right)^2\right] = \frac{1}{T-1} E_b\left[\left(\sum_{K=1}^{T-1} \prod_{t=0}^{K} \frac{\pi}{b} G_0\right)^2\right] \rightarrow \infty$

$\frac{1}{T-1}$ constant,

because $E_b\left[\left(\sum_{K=1}^{T-1} \prod_{t=0}^{K} \frac{\pi}{b} G_0\right)^2\right] \rightarrow \infty$
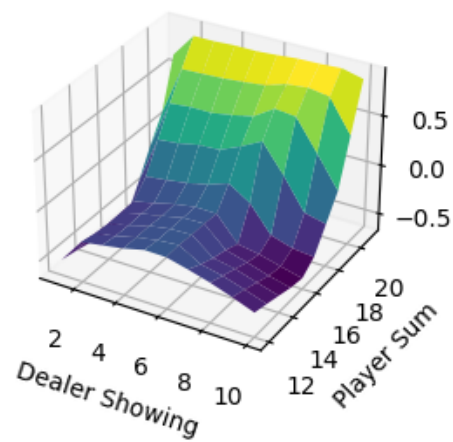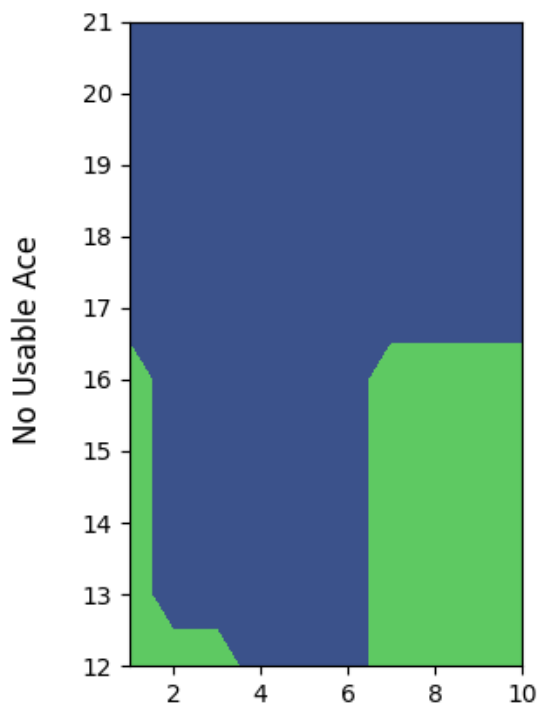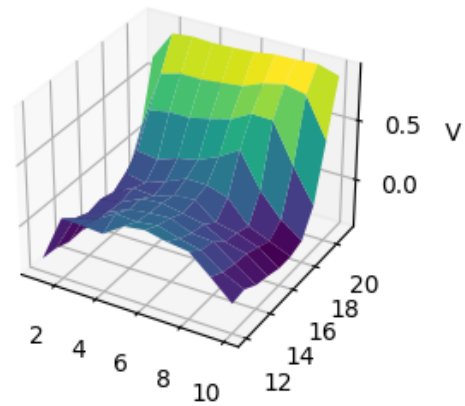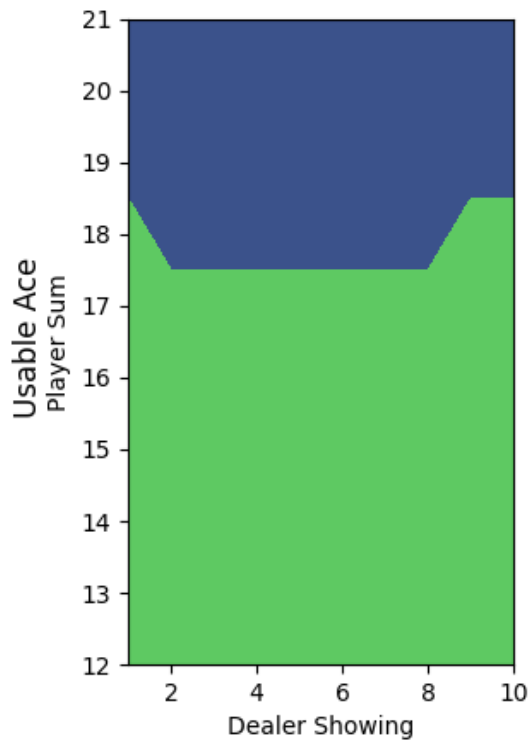
and $\bar{X}^2$ for every-visit is finit too

**For running Q3 and Q4, in algorithm.py, there is a main function in which you need to uncomment each question and part you like to run**

# Q3- a



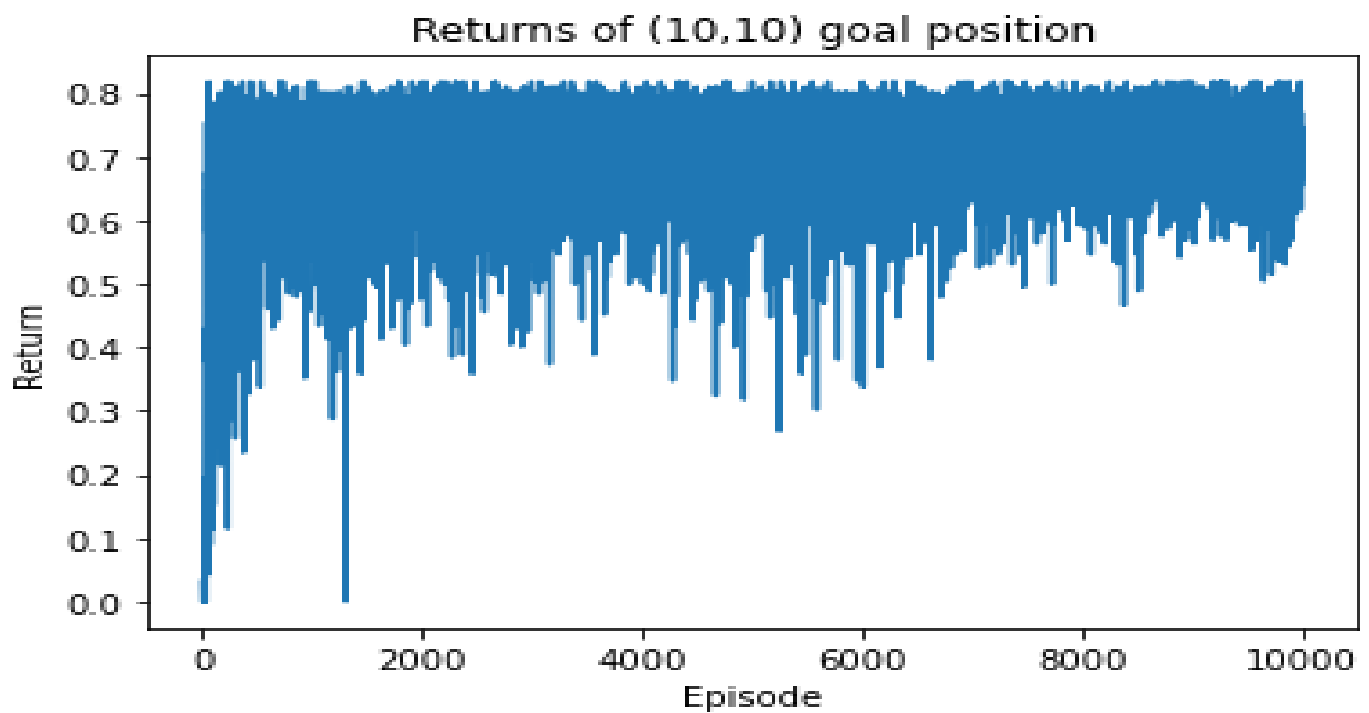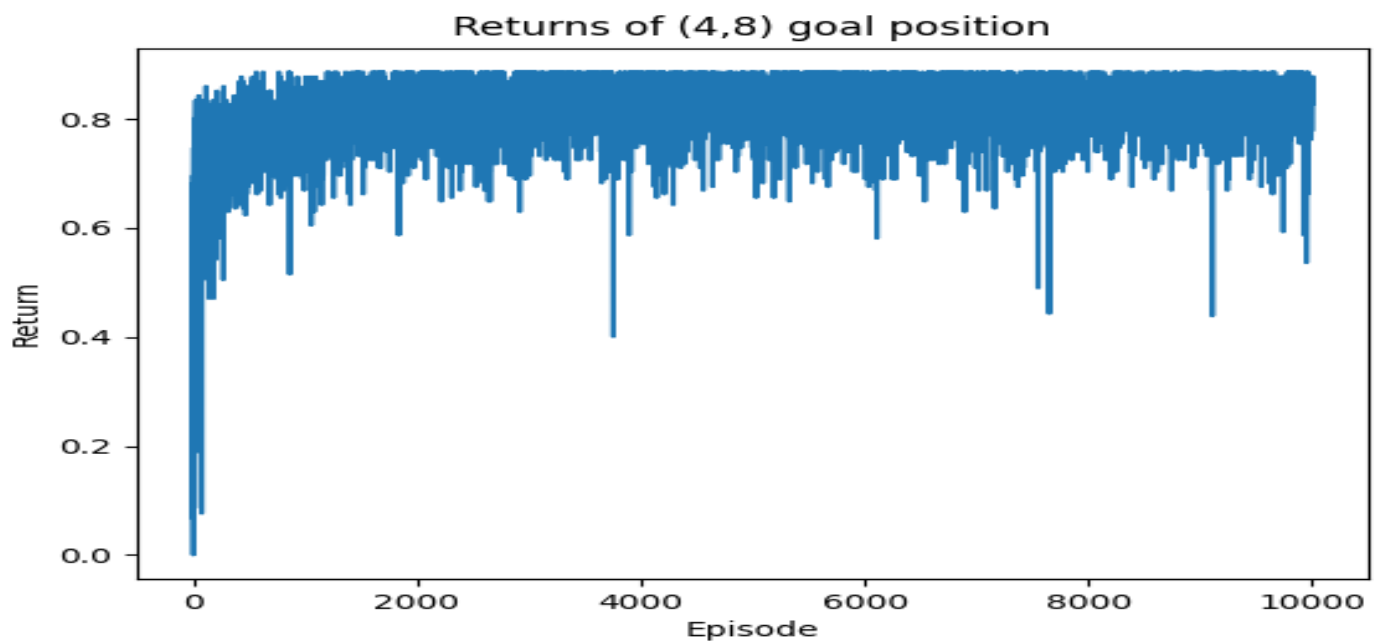After 10,000 and After 500,000 Episodes

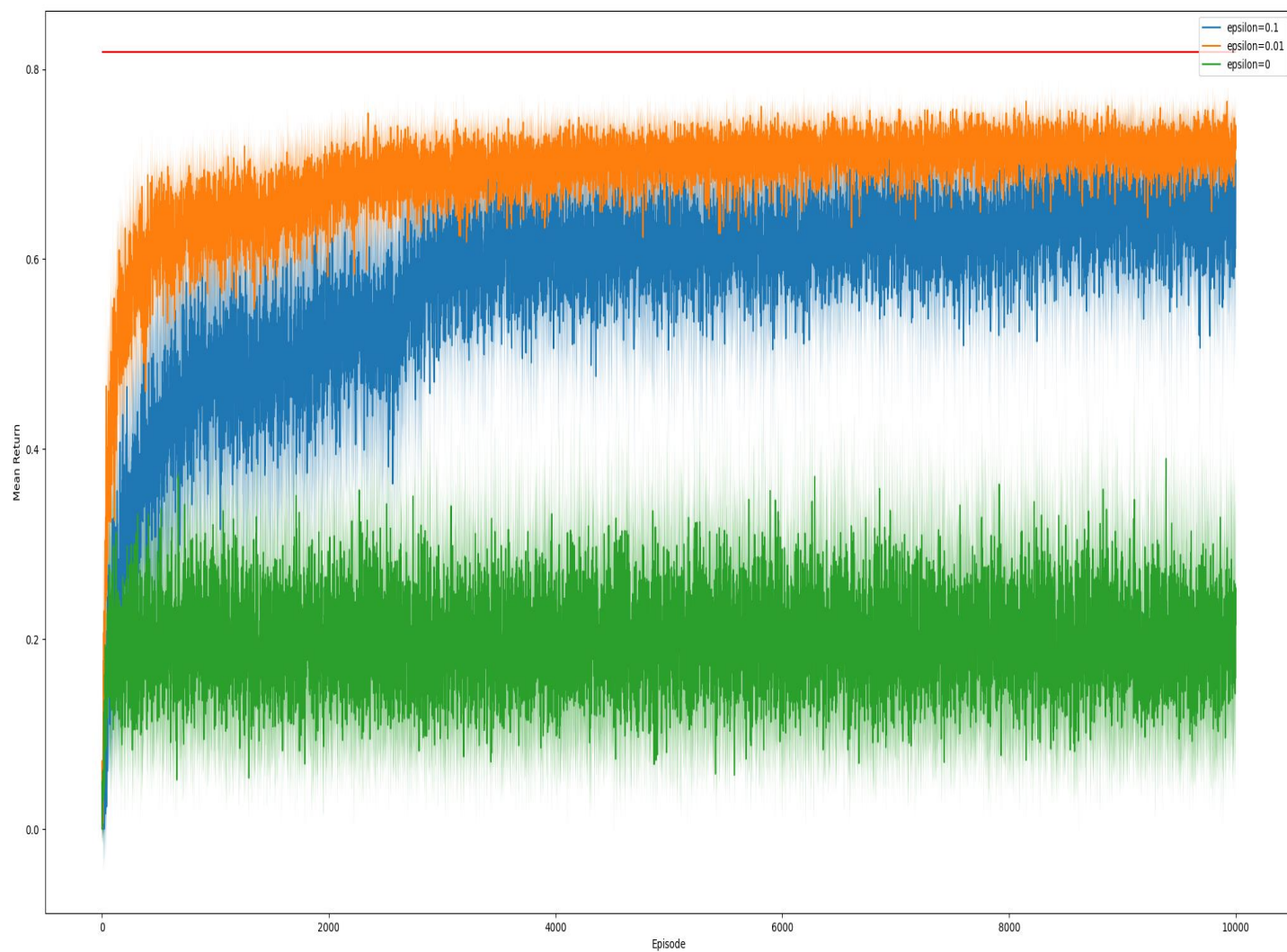# Q3- b



Green represents hit and blue represents stay.

# Q4- a -

**The following plot shows returns when the goal is at (4,8), showing that it works for non-original goal states**



Returns of (4,8) goal position
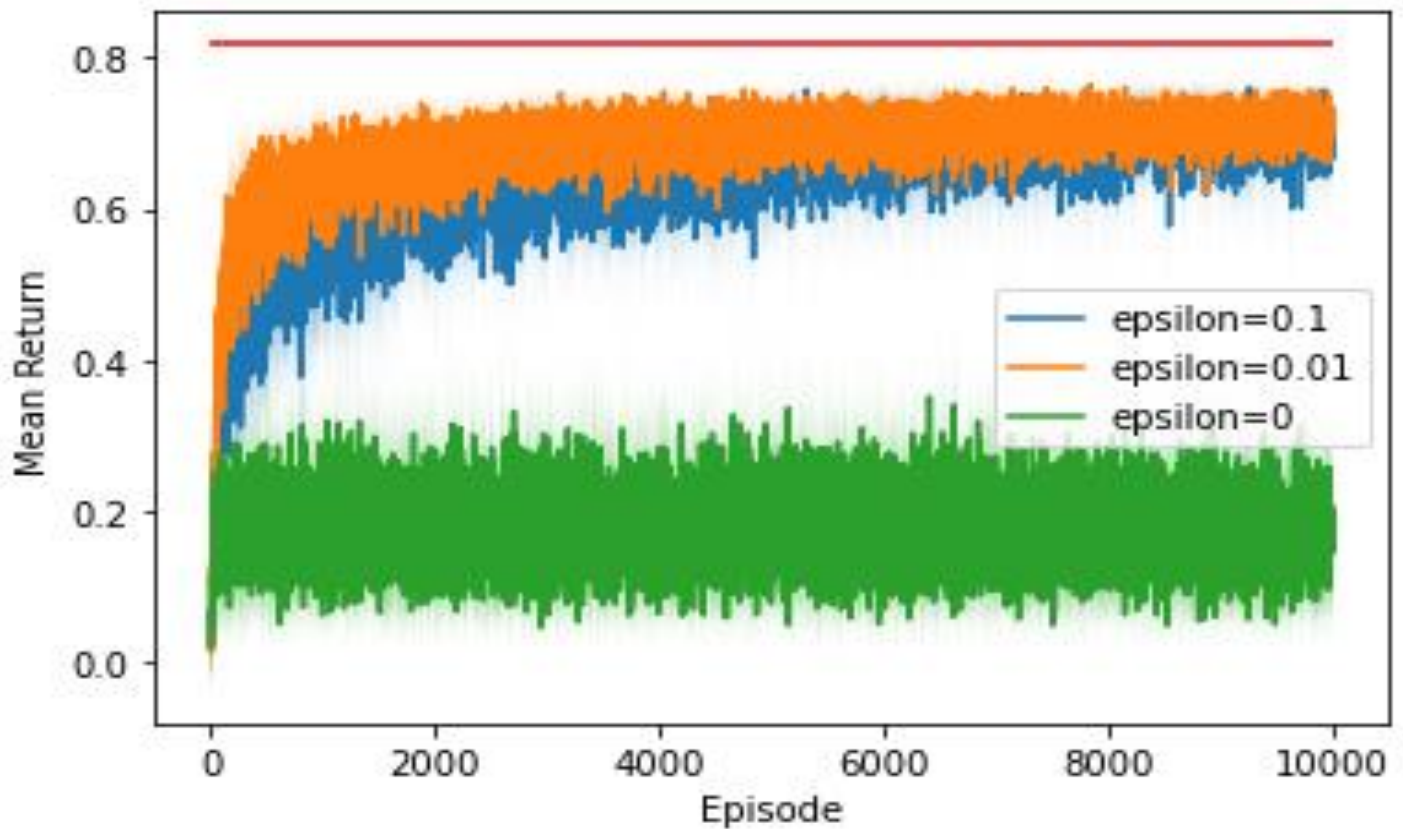


Returns of (10,10) goal position

**b-**

Learning curve of (4,8) goal position

Learning curve of (10,10) goal position

C) Without exploring starts, the agent will follow the path it first finds to the goal as the only exploration. Because the initial starter is random, there is a large chance that this path will be very inefficient, far from optimal answer, which leads to the mean return of only 0.2

5-a) $V_{n+1} = \dfrac{\sum^n W_K G_K}{\sum^n W_K} = \dfrac{W_n G_n + \sum^{n-1} W_K G_K}{\sum^n W_K} \cdot \dfrac{\sum^{n-1} W_K}{\sum^{n-1} W_K}$

$= \left[\dfrac{W_n G_n}{C_{n-1}} + V_n\right] \dfrac{C_{n-1}}{C_n} = \dfrac{W_n G_n}{C_n} + \dfrac{V_n C_{n-1}}{C_n} \qquad C_i = \sum_{K=1}^{i} W_K$

$= V_n + \dfrac{W_n G_n}{C_n} + \dfrac{V_n C_{n-1}}{C_n} - V_n = V_n + \dfrac{W_n G_n}{C_n} + \dfrac{V C_{n-1} - V_n C_n}{C_n}$

$= V_n + \dfrac{W_n G_n}{C_n} + \dfrac{-V_n W_n}{C_n} = V_n + \dfrac{W_n}{C_n}\left[G_n - V_n\right]$

b) because we are assuming that the target
policy is greedy and deterministic so we
could consider $\pi(A_t | s_t) = 1$

$\Rightarrow \dfrac{\pi}{b} = \dfrac{1}{b}$

**Q6)** I am still working on question 6, when ever it finish I will send it.