

# RL – Reinforcement Learning

## CSE564

### Homework 1

Name: Ansh Arora

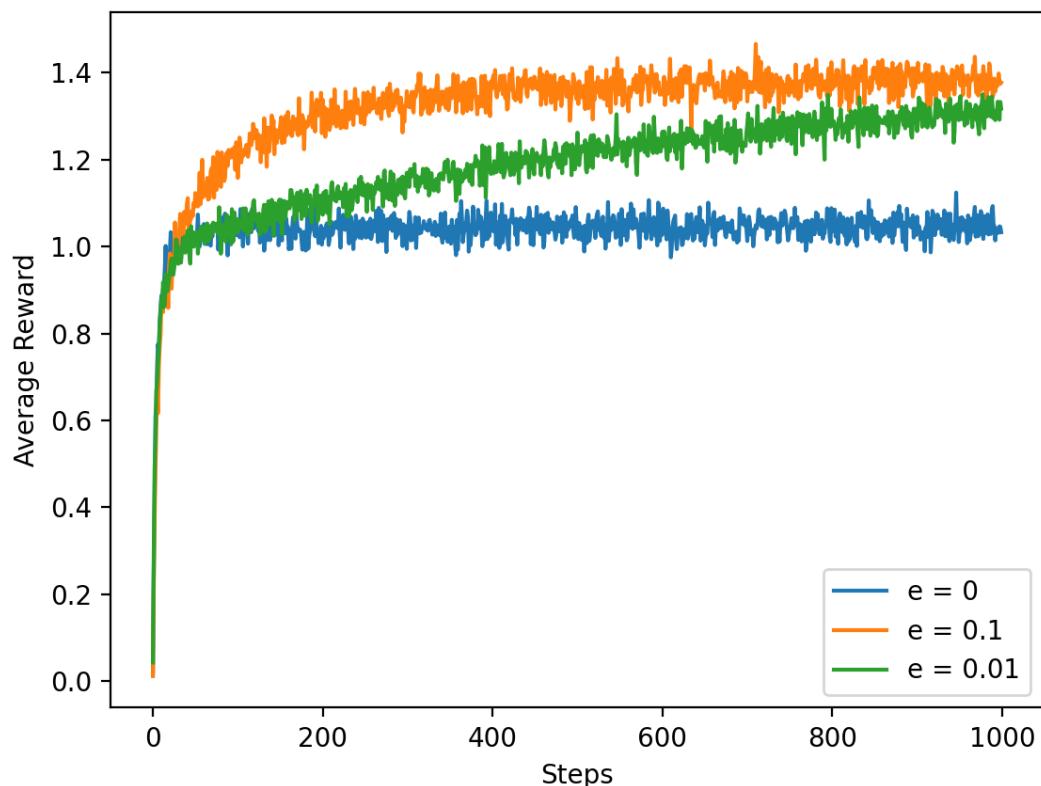
Roll No: 2019022

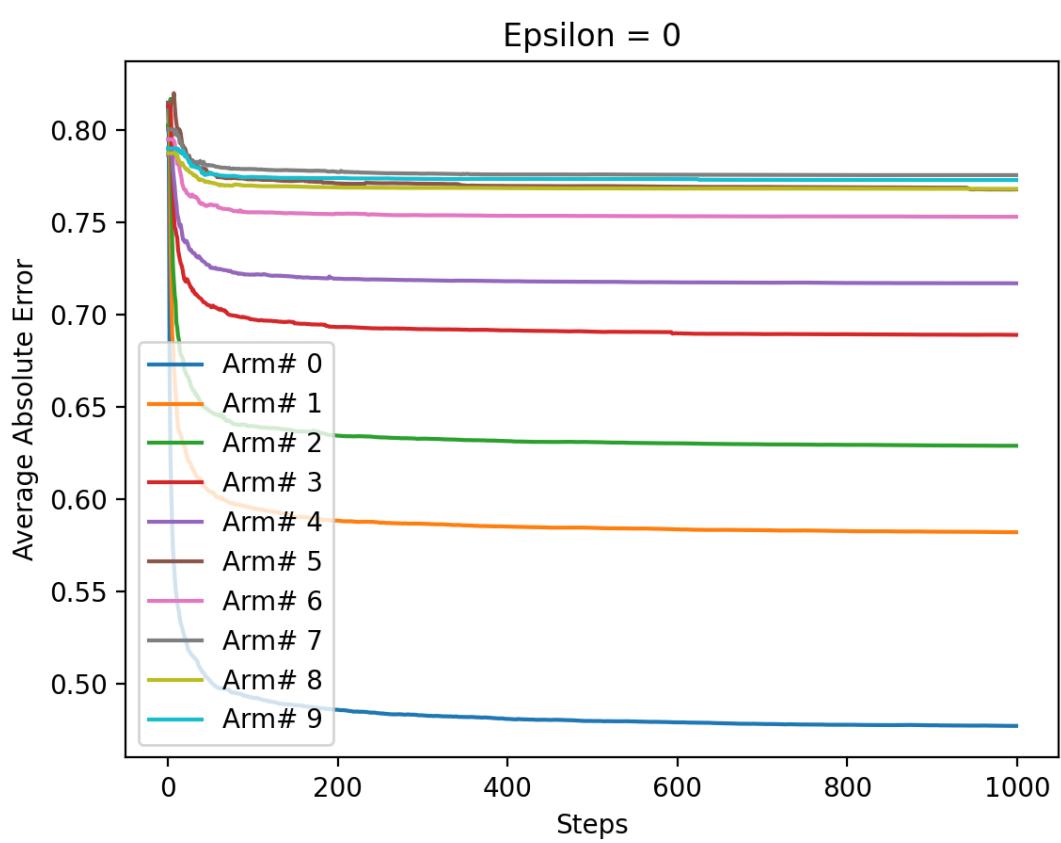
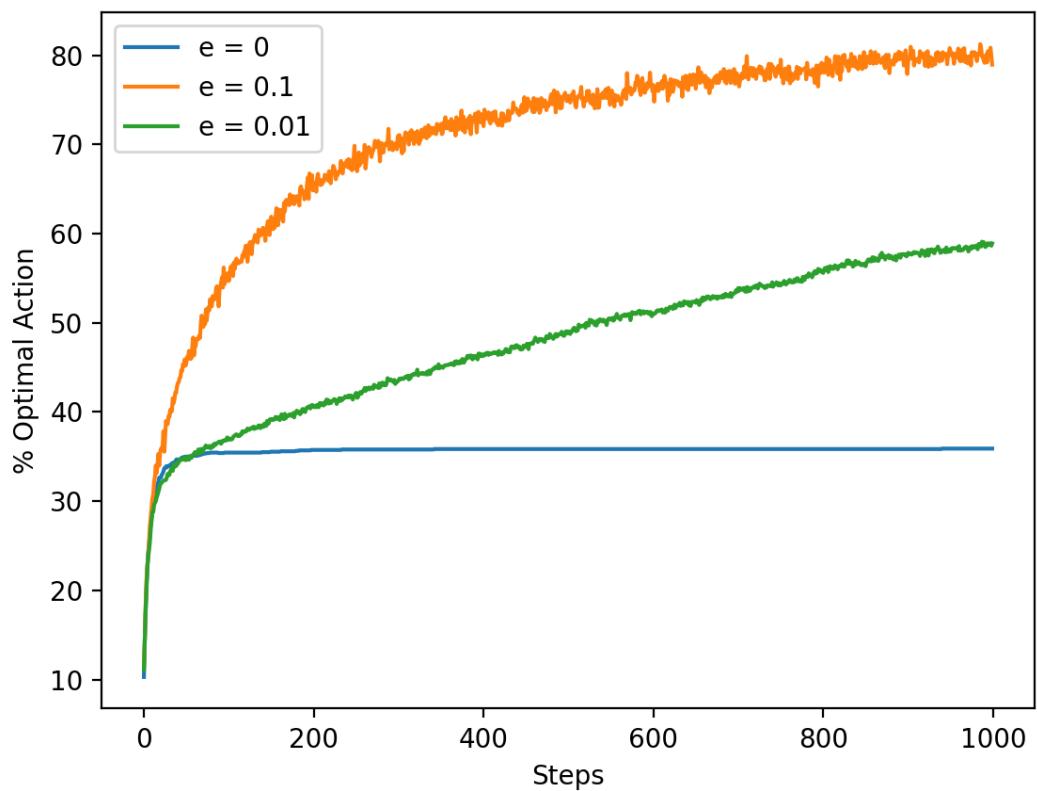
Graphs pasted at beginning, theory Qs at end.

#### Question 1

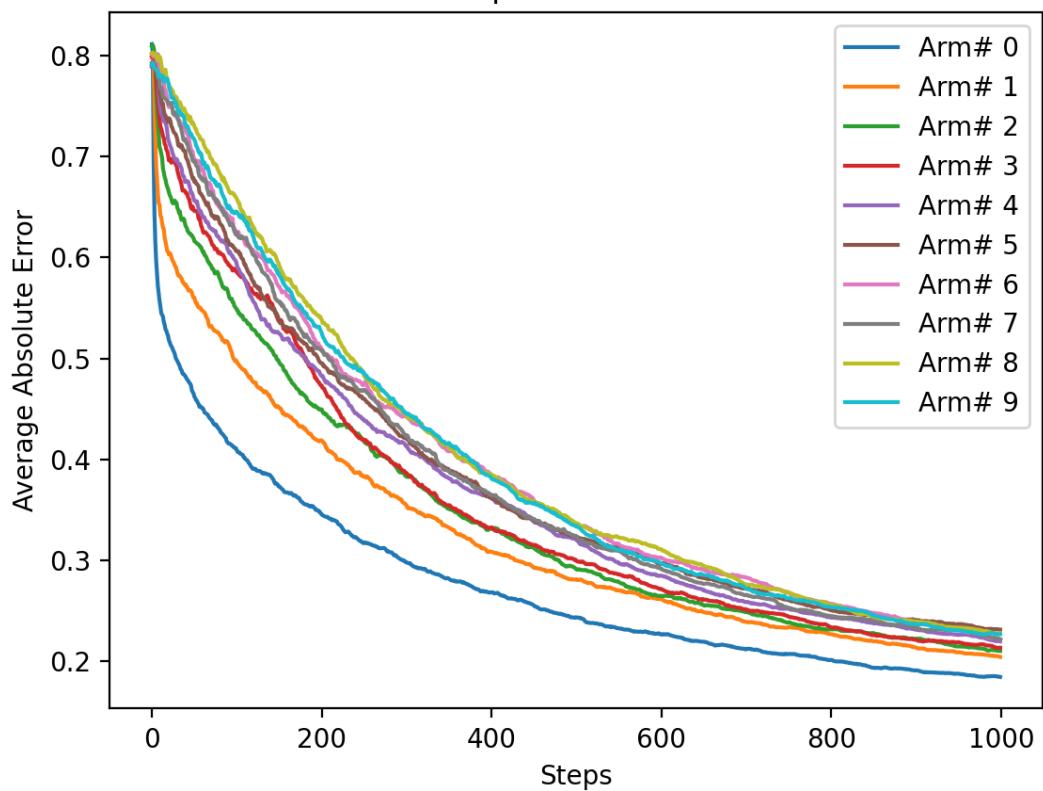
For epsilon = 0, 0.1, 0.01

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q1/Q1.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q1/Q1.py)

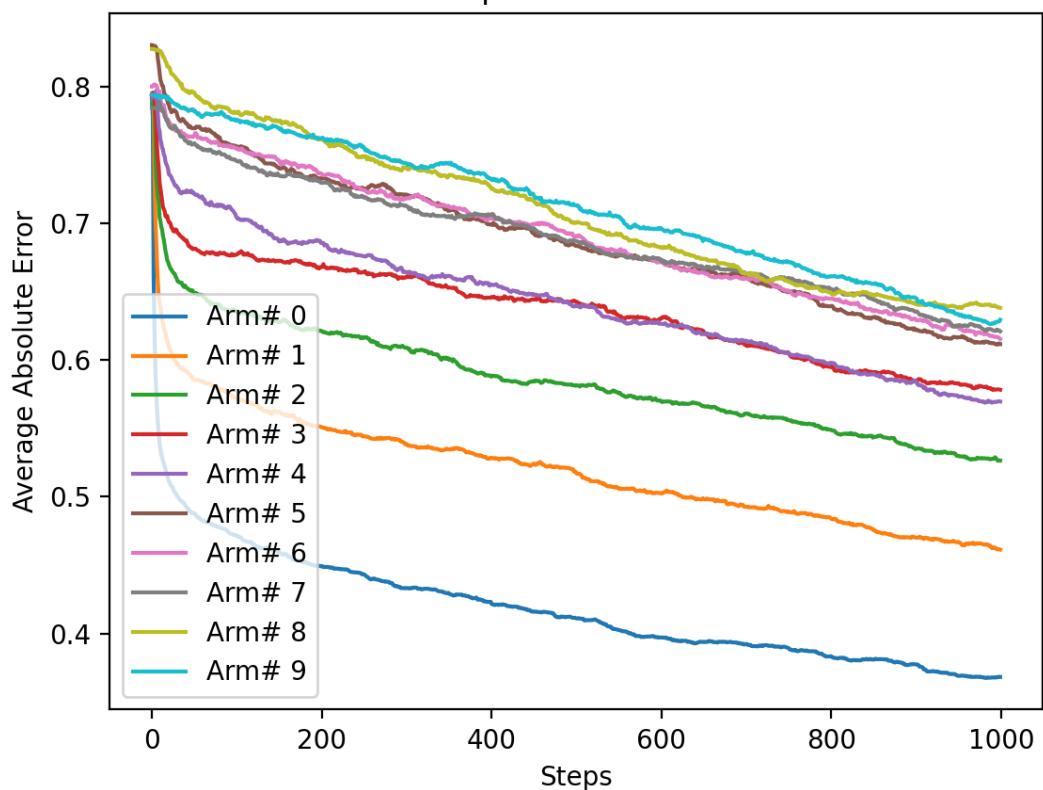




Epsilon = 0.1

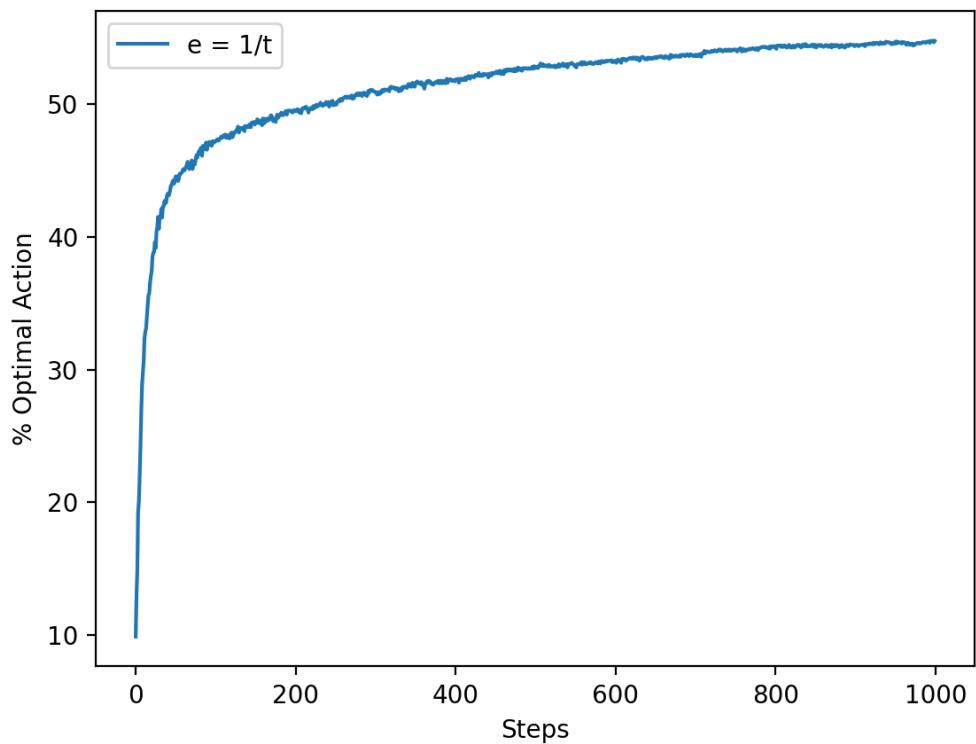
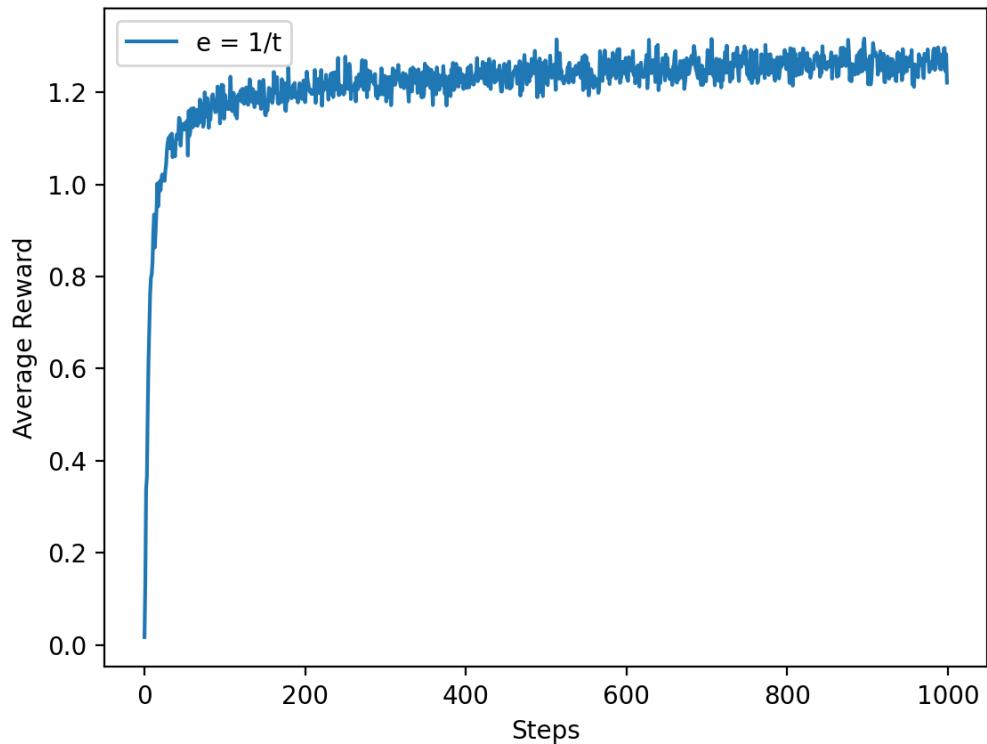


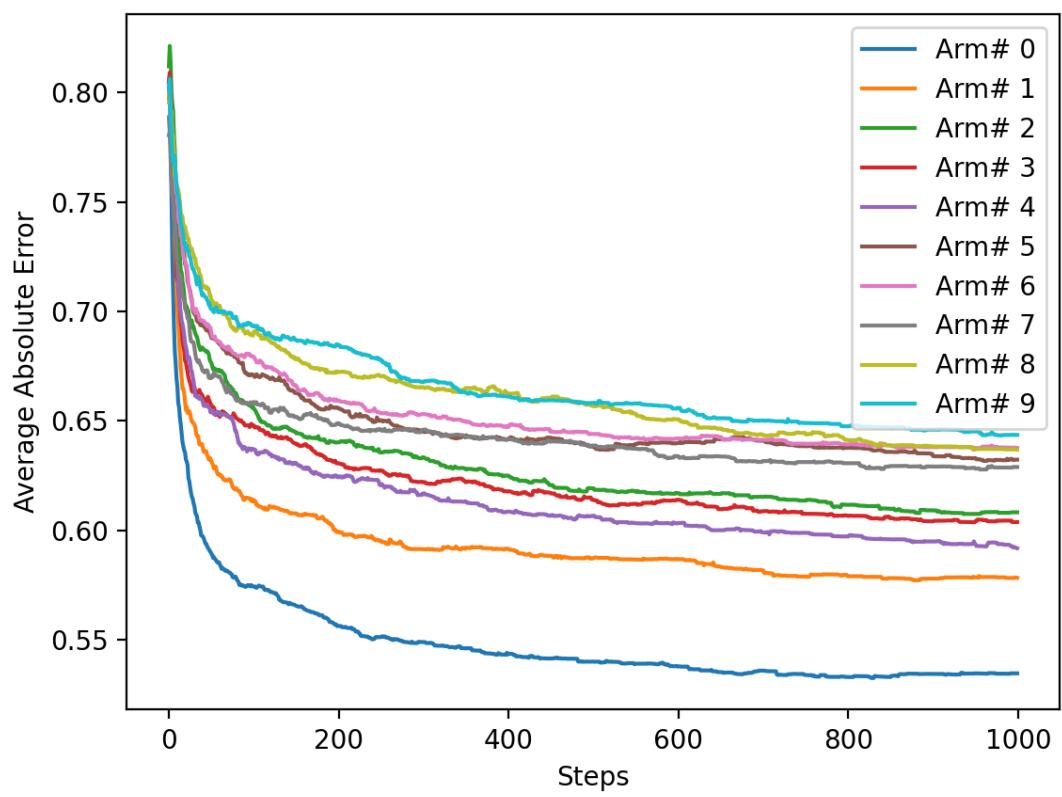
Epsilon = 0.01



For epsilon = 1/t

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q1/Q1b.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q1/Q1b.py)

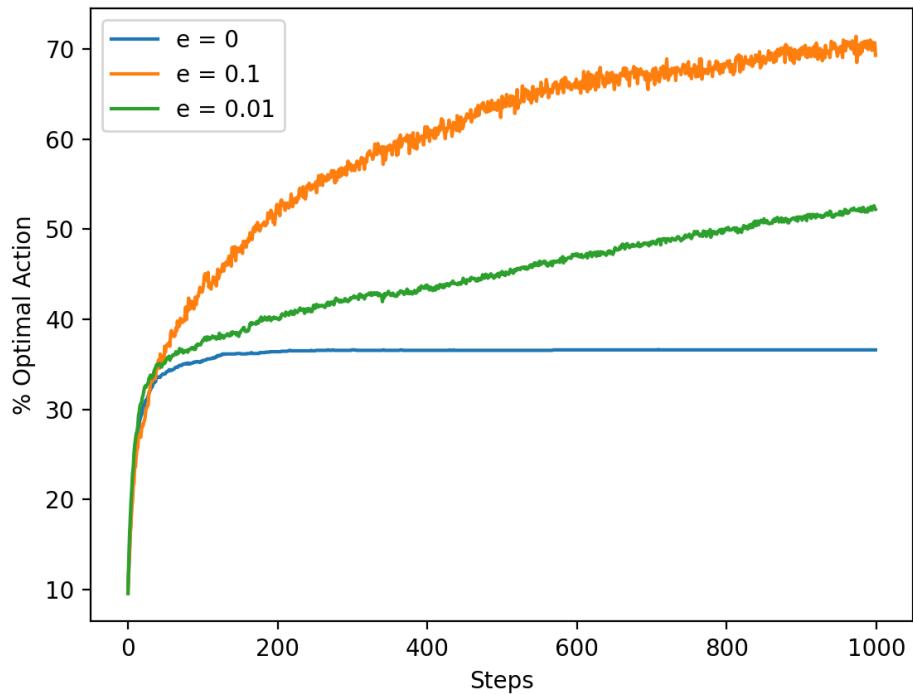
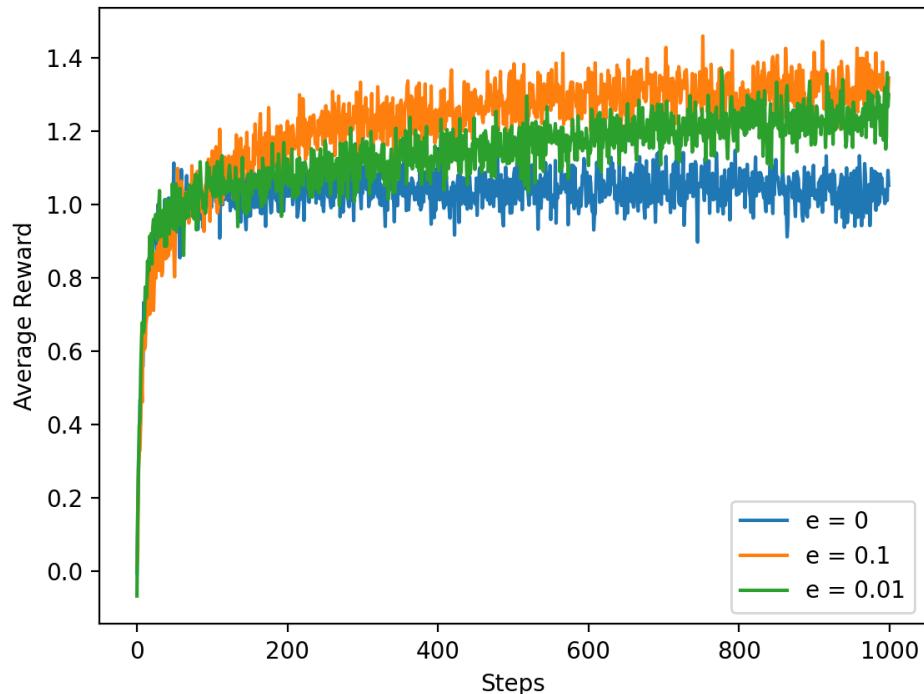




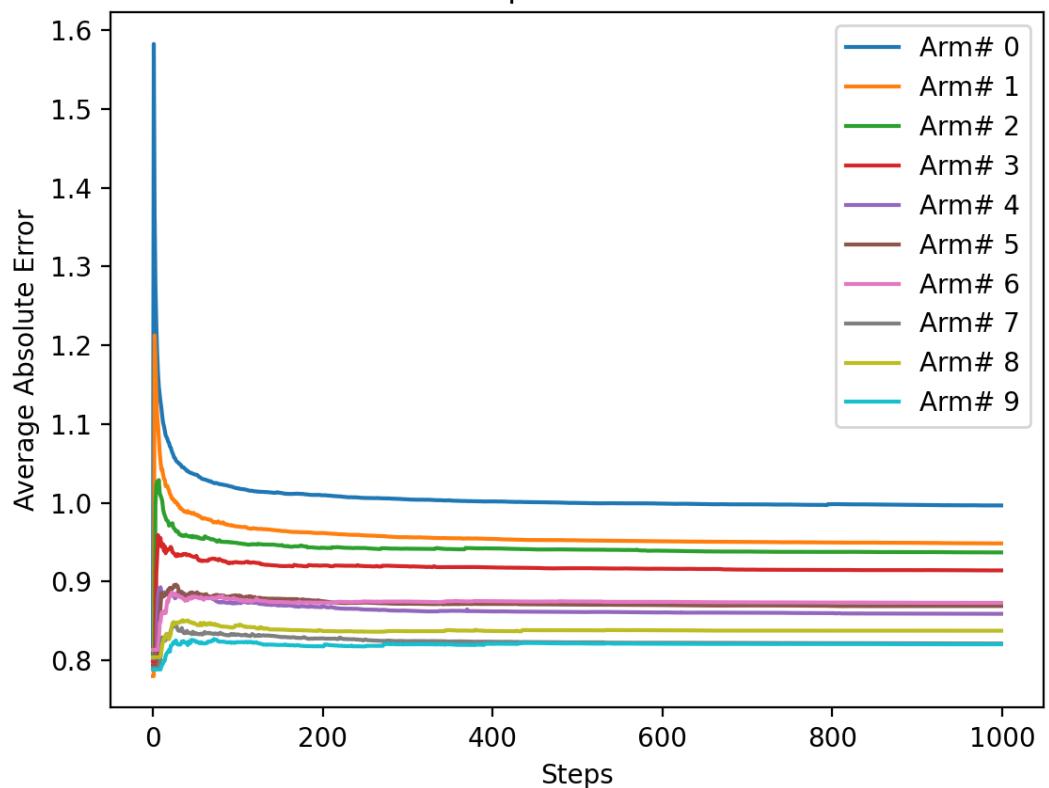
## Question 2

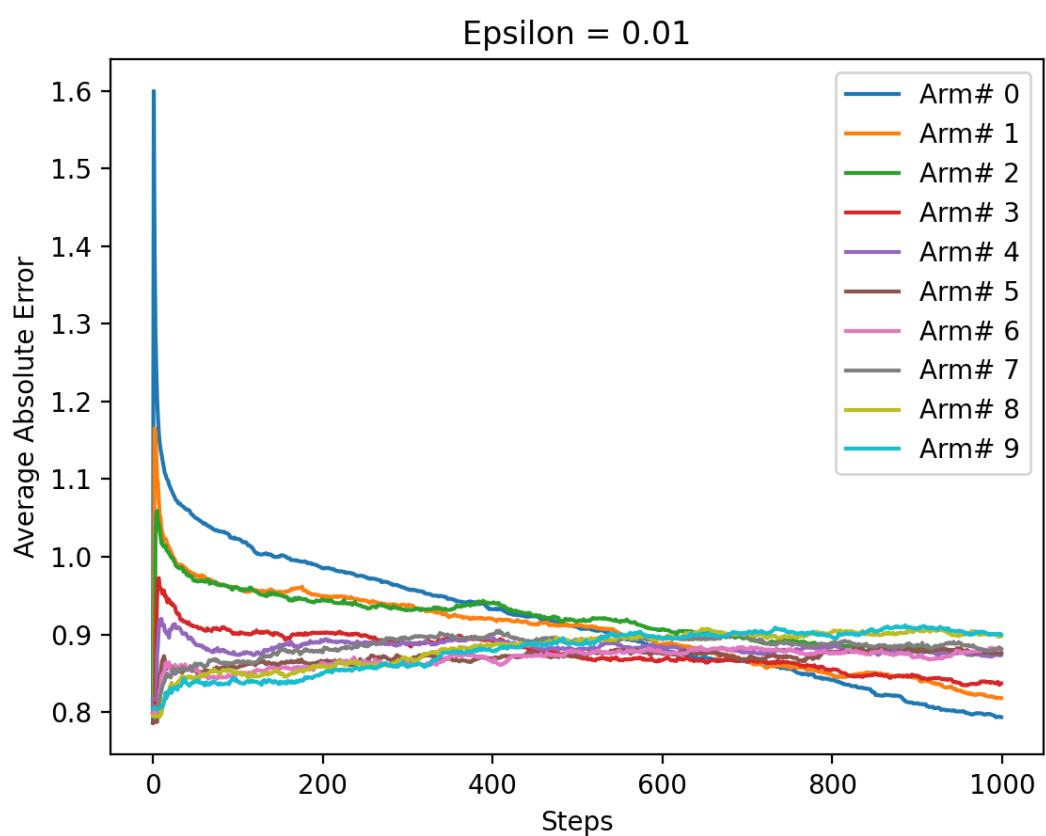
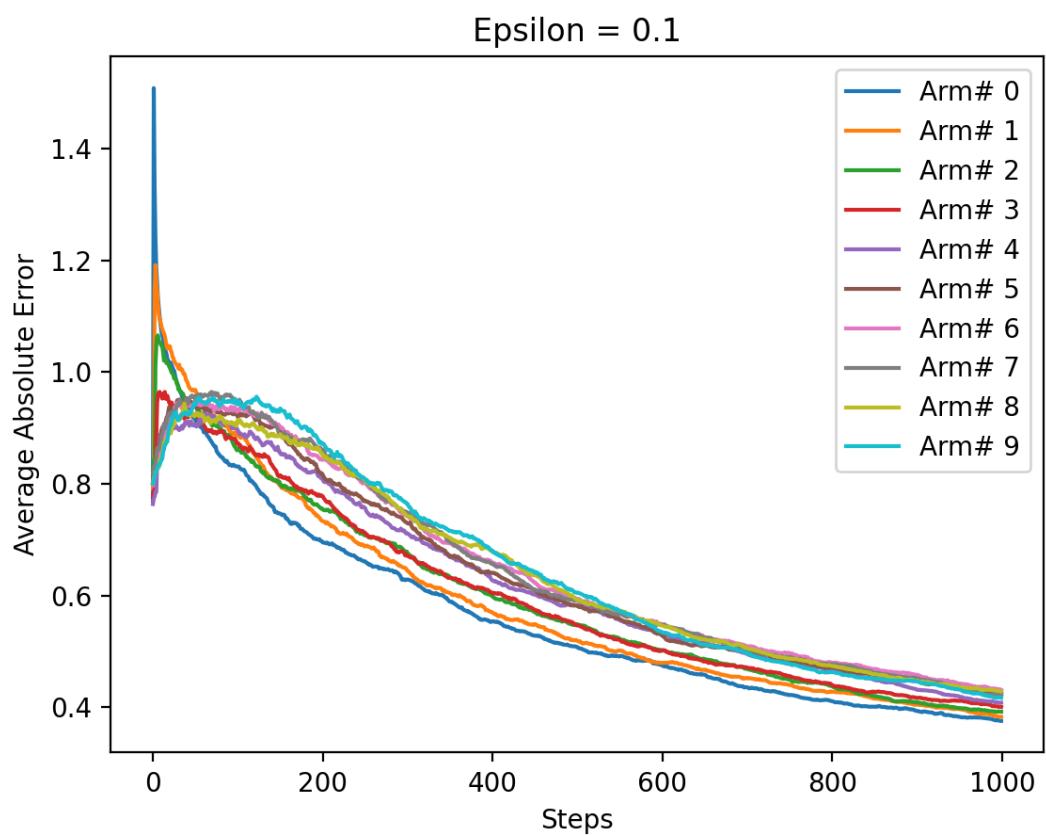
For epsilon = 0, 0.1, 0.01

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q2/Q2.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q2/Q2.py)



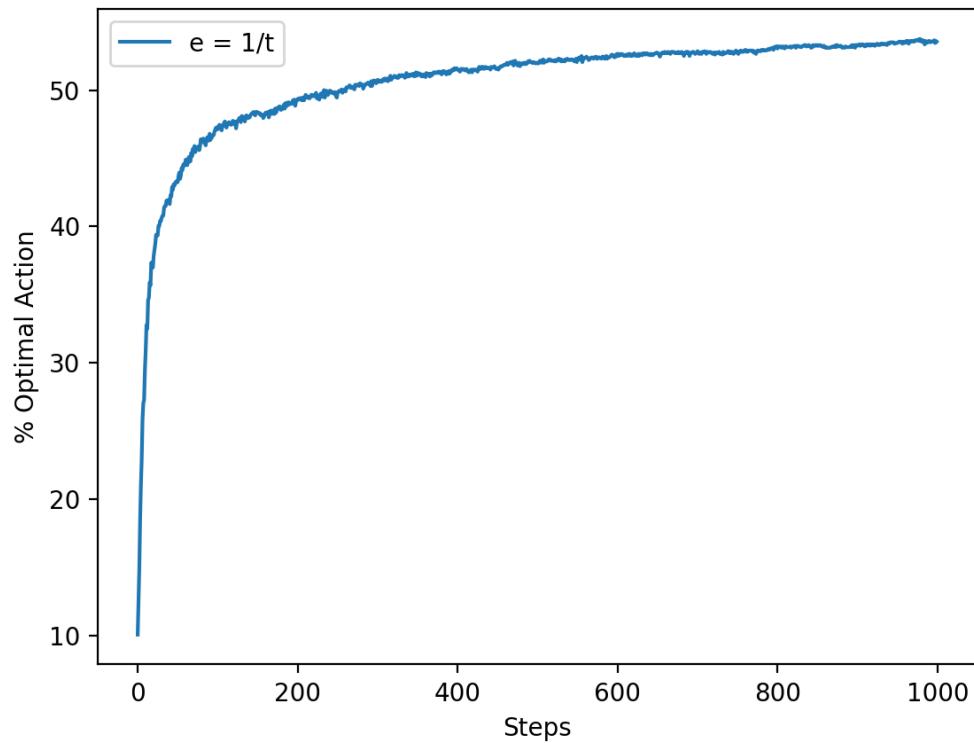
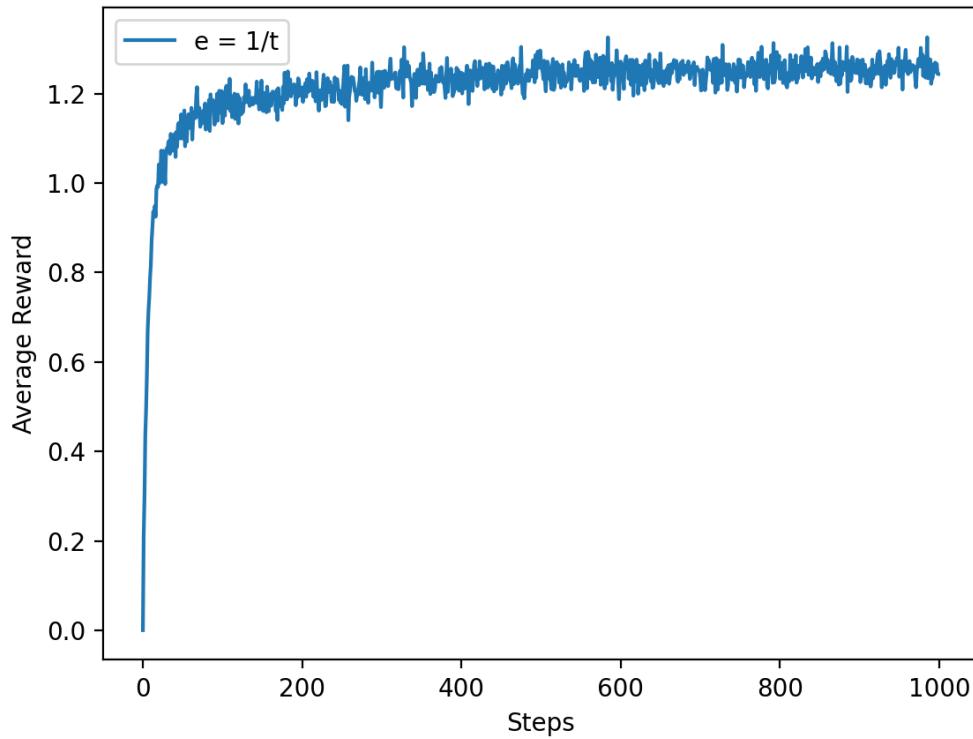
Epsilon = 0

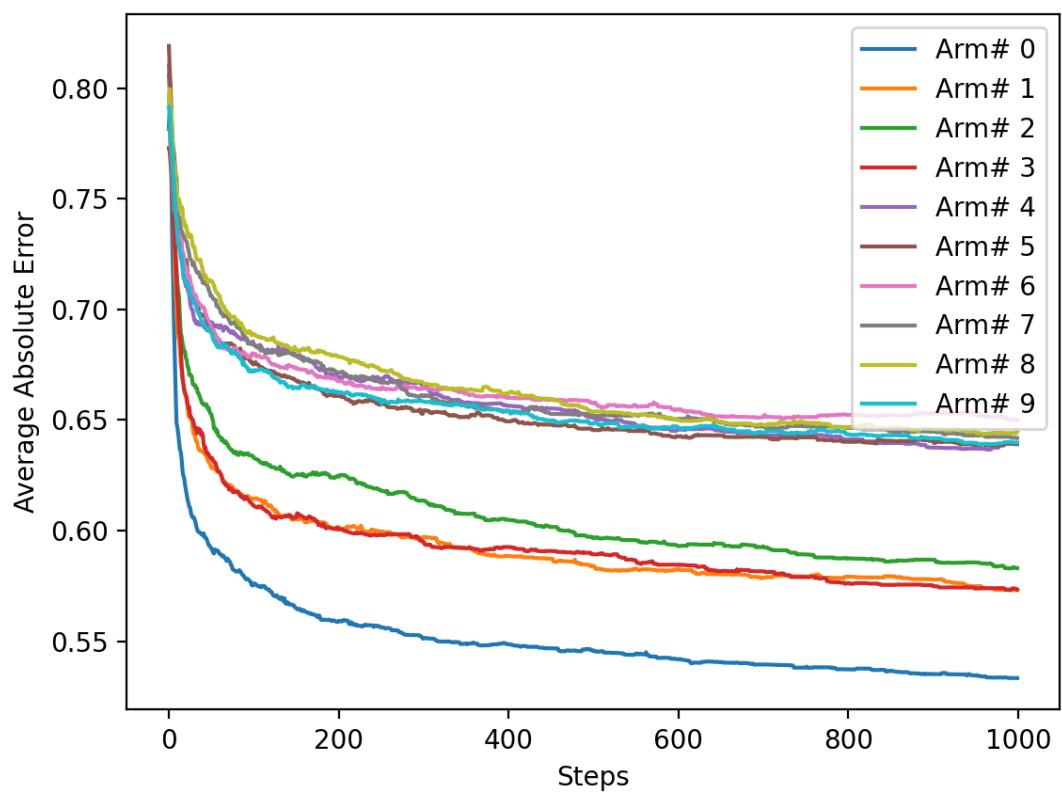




For epsilon =  $1/t$  -

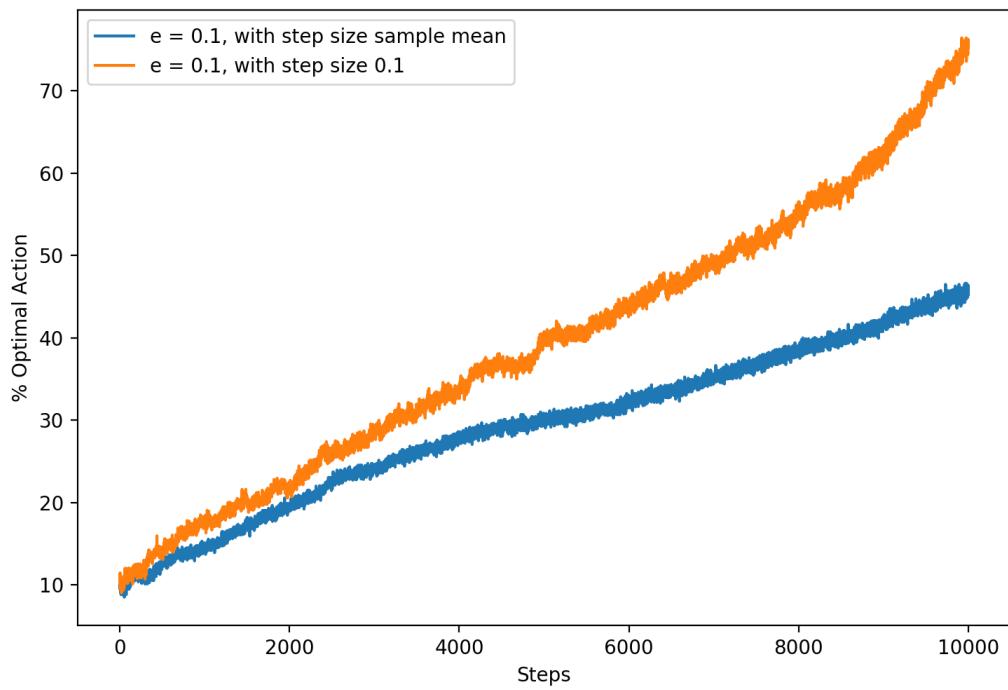
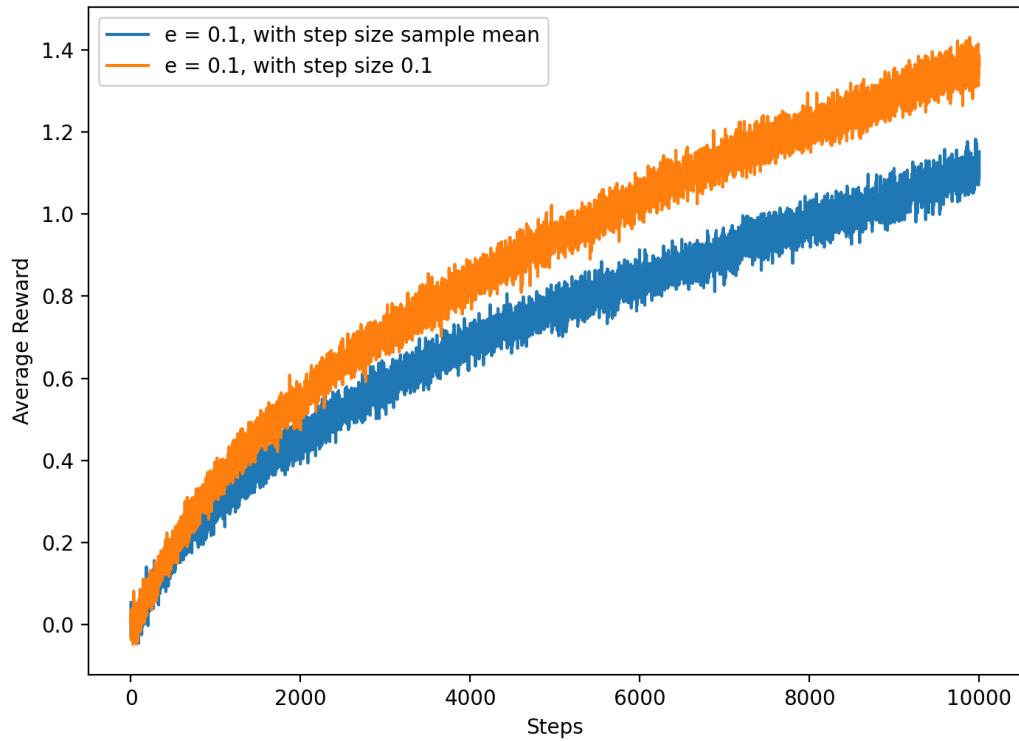
[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q2/Q2b.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q2/Q2b.py)





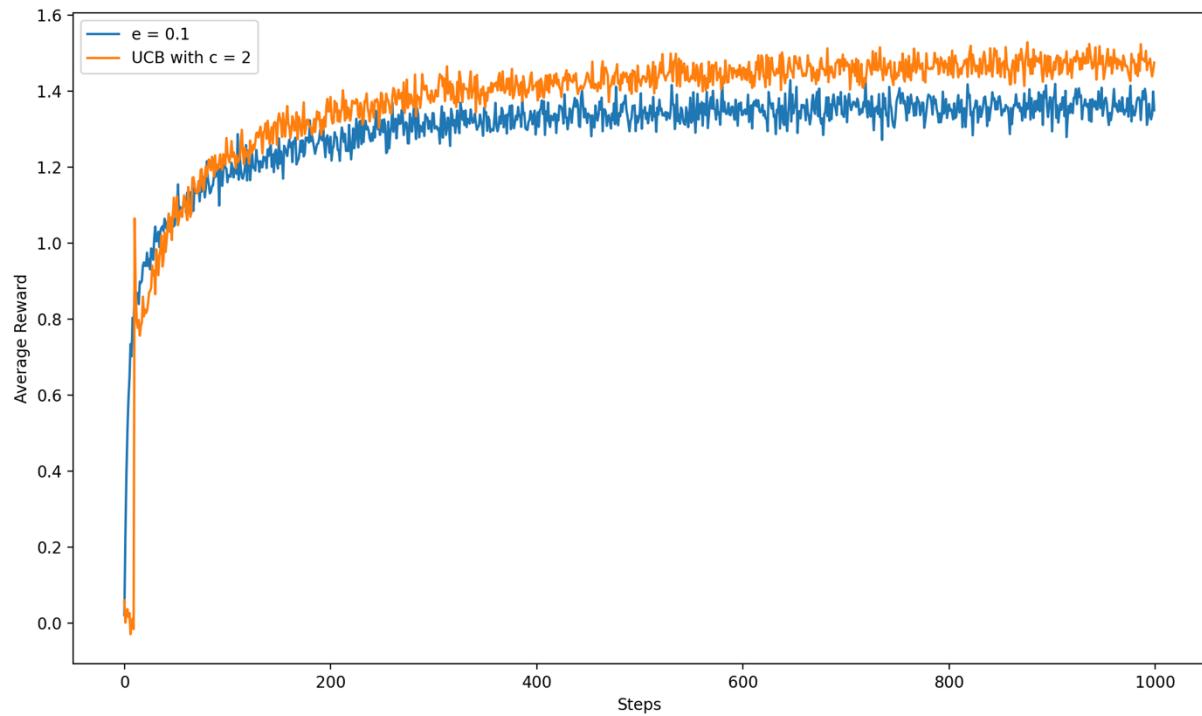
## Question 5 –

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q5/Q5.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q5/Q5.py)



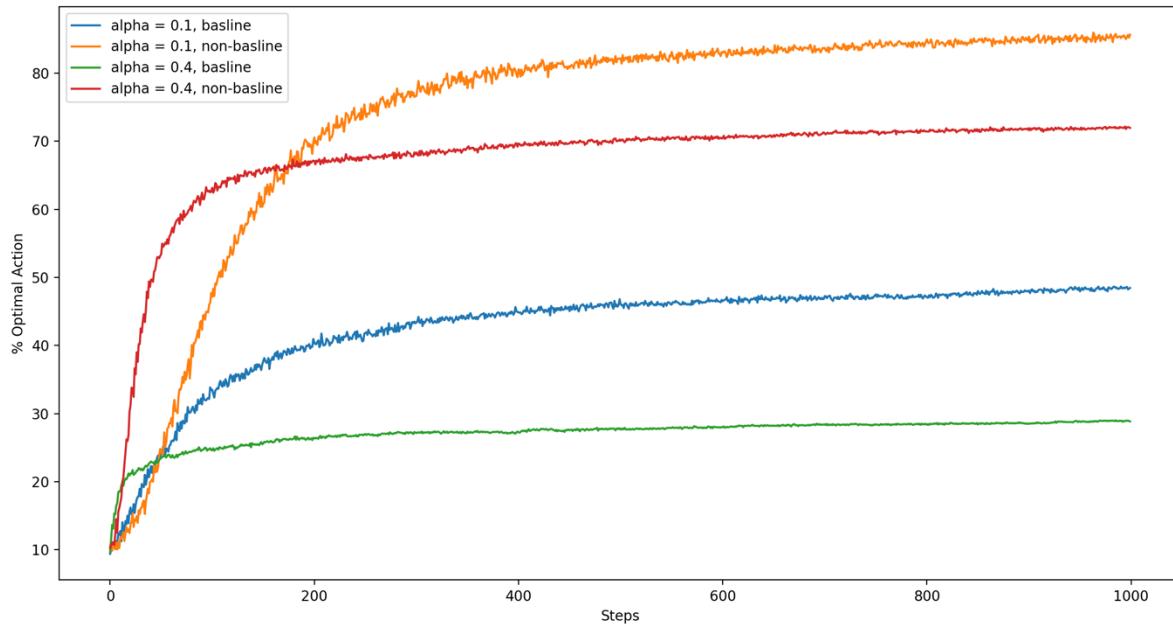
## Question 6 -

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q6/Q6.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q6/Q6.py)



## Question 7 -

[https://github.com/arora-ansh/RL-M2021\\_Ansh\\_2019022/blob/main/HW1/Q7/q7.py](https://github.com/arora-ansh/RL-M2021_Ansh_2019022/blob/main/HW1/Q7/q7.py)



# RL - Reinforcement Learning

Ansh Arora

2019022

Date: \_\_\_/\_\_\_/\_\_\_

- ③ Which method amongst  $\epsilon = 0.1$ ,  $\epsilon = 0.01$ ,  $\epsilon = 0$  and  $\epsilon = 1/t$  will perform the best.

For best performing, Expected value of that particular run should be highest. This would be highest when ~~the~~ probability of picking the most optimal arm at infinity would be high. Let this be  $P(x)$ .

i)  $\epsilon = 0.1$

$$\begin{aligned} P(x) &= 1 - \epsilon + \frac{\epsilon}{K \text{ (no. of arms)}} \\ &= 1 - 0.1 + \frac{0.1}{10} \\ &= 0.9 + 0.01 = 0.901 \end{aligned}$$

ii)  $\epsilon = 0.01$

$$\begin{aligned} P(x) &= 1 - \epsilon + \frac{\epsilon}{K} \\ &= 1 - 0.01 + \frac{0.01}{10} \\ &= 0.99 + 0.001 \\ &= 0.991 \end{aligned}$$

iii)  $\epsilon = 0$

$$P(x) = 1$$

However, in this particular case since there will be no exploration, there is a  $1/10$  chance that we would be able to pick the right arm greedily making this inefficient.

7)  $\varepsilon = \cancel{\frac{1}{t}} - \frac{1}{t}$

$$P(x) = \lim_{n \rightarrow \infty} 1 - \frac{1}{t} + \frac{1}{t^K}$$

$$= 1$$

$\therefore \varepsilon = 1/t$  would be the best performing method in the long run in terms of cumulative reward.

- ④ In case of sample mean after the action has been chosen once the dependence on initial choice goes away. This can be shown as follows -

~~$$Q_{t+1}(a) = Q_t(a) + \frac{1}{t} [R_t(a) - Q_t(a)]$$~~

$$\text{At } t=1$$

$$Q_2(a) = Q_1(a) + \cancel{\frac{1}{1}} [R_1(a) - Q_1(a)]$$

$$\Rightarrow Q_2(a) = R_1(a)$$

Now,

$$\text{At } t=2$$

$$\begin{aligned} Q_3(a) &= Q_2(a) + \frac{1}{2} [R_2(a) - Q_2(a)] \\ &= \frac{Q_2(a)}{2} + \frac{R_2(a)}{2} \end{aligned}$$

and so on.

since all future terms would now be dependent on  $Q_1(a)$ , and  $Q_2$  is independent of  $Q_1(a)$  sample mean is not influenced by initial choice of  $Q_1(a) + a$ .

Now taking a constant step size  $\alpha$ , we have

$$\begin{aligned}
 Q_{t+1}(a) &= Q_t(a) + \alpha [R_t(a) - Q_t(a)] \\
 &= \alpha R_t(a) + (1-\alpha) Q_t(a) \\
 &= \alpha R_t(a) + (1-\alpha)[\alpha R_{t-1}(a) + (1-\alpha) Q_{t-1}(a)] \\
 &= \alpha R_t(a) + (1-\alpha)\alpha R_{t-1}(a) + (1-\alpha)^2 Q_{t-1}(a) \\
 &= (1-\alpha)^n Q_1(a) + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i
 \end{aligned}$$

(2.6  
from book)

Thus as can be seen above it will be dependent on  $Q_1(a)$ , the initial choice  $a$ .

For  $\alpha \in (0, 1)$ , if  $\alpha$  value is reduced  $(1-\alpha)$  increases, as a result of which the coefficient of initial choice  $Q_1(a)$   $(1-\alpha)^n$  increases. Thus, the dependence is stronger for a smaller  $\alpha$ .

To deal with dependence on  $Q_1(a)$ , we can select the constant value to be 1.

This would mean  $(1-\alpha)^n = 0^n = 0$  would be coefficient of  $Q_1(a)$ , thus making  $Q_{t+1}(a)$  independent of  $Q_1(a) + a$ .

⑥ In the Upper Confidence Bound method,

$$A_t = \operatorname{argmax}_a [Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}}]$$

Here  $N_t(a)$  denotes the number of times that action  $a$  gets selected before time  $t$ . At  $N_t(a) = 0$ ,  $a$  is considered to be a maximizing action.

In the first 10 steps, on an average the agent would cycle through all these actions (since it is taken to be maximal). On the 11th action, the agent will end up choosing greedily, leading to a sudden spike value.

Now, on the 12th step,  $\sqrt{\frac{\ln t}{N_t(a)}}$  is equal to  $\sqrt{\frac{\ln 12}{10}}$ , for 9 of the 10 actions, while it is  $\sqrt{\frac{\ln 12}{1}}$  for 1 of the actions.

that was picked greedily on the 11th step. Since this value is  $\frac{1}{\sqrt{2}}$  times the last value

that was picked on step 11, we can say that on average the pick would be lesser than step 11, hence leading to a fall in value.

∴ A spike is obtained on Step 11.