

Reinforcement Learning

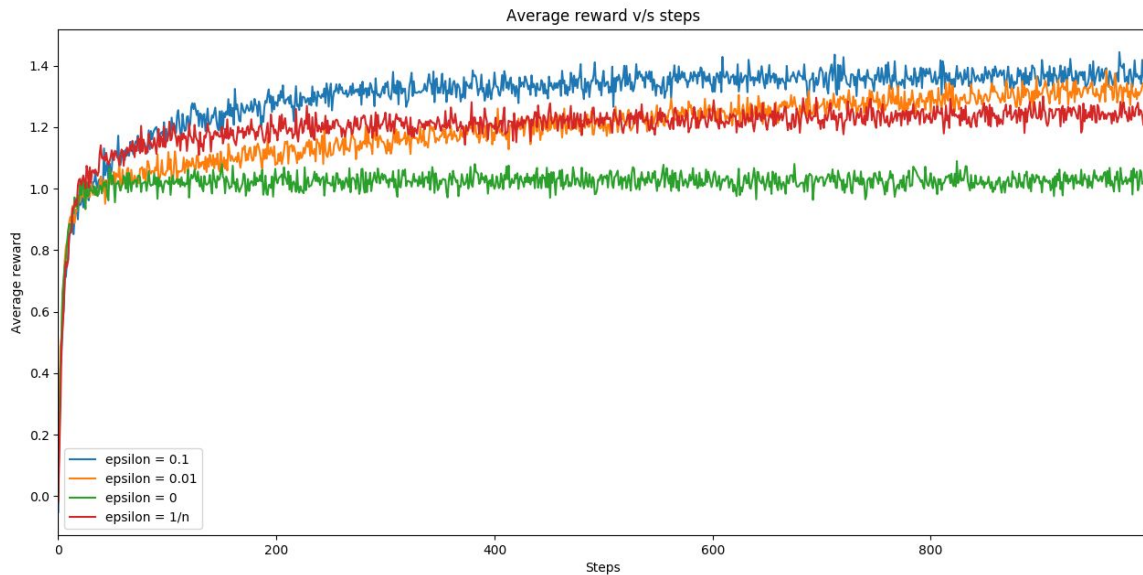
Assignment 1

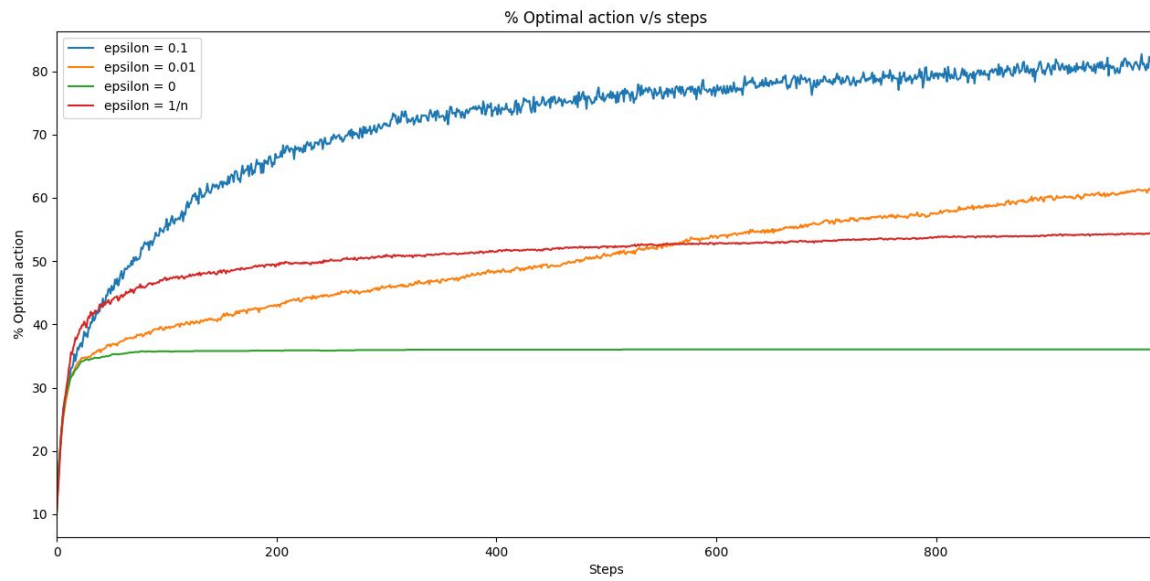
Setu Gupta (2018190)

Q1.)

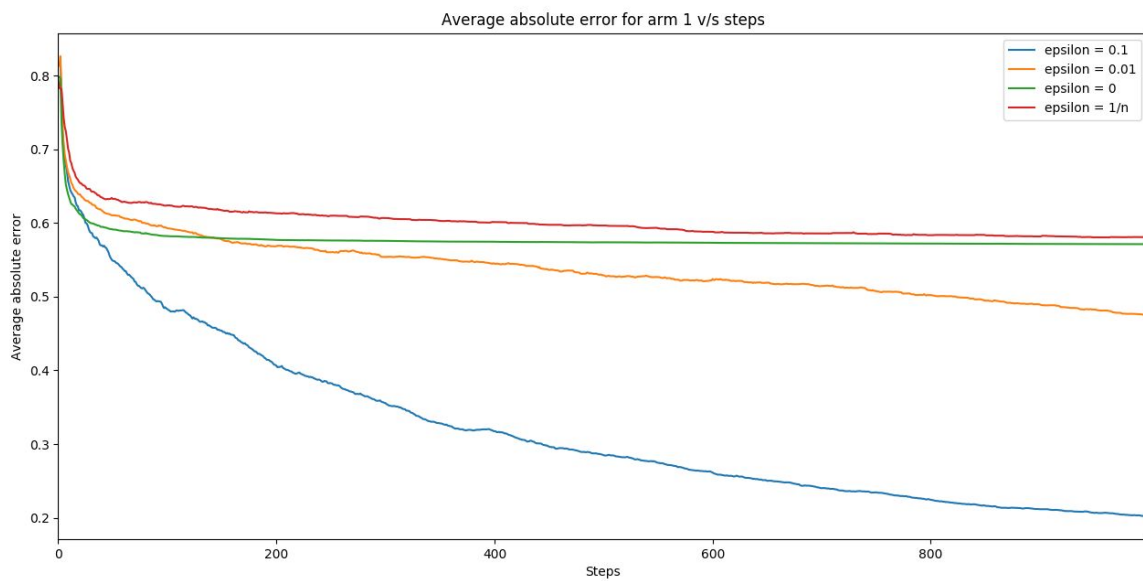
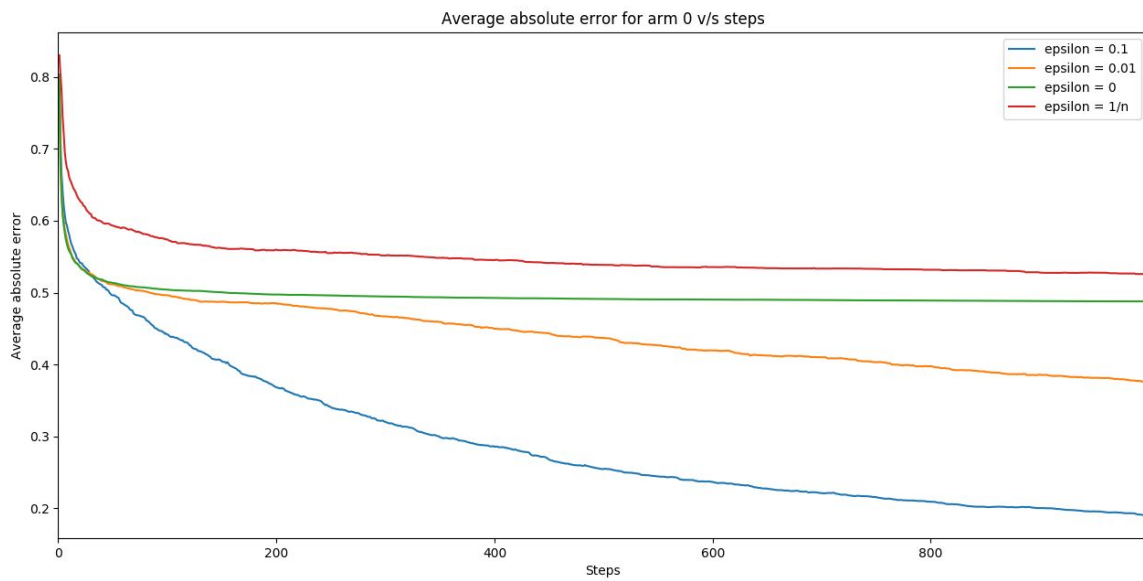
$\epsilon = 1/n$ is a sequence of epsilon which satisfies Eq 2.7 of SB. Here n is the number of current timestep.

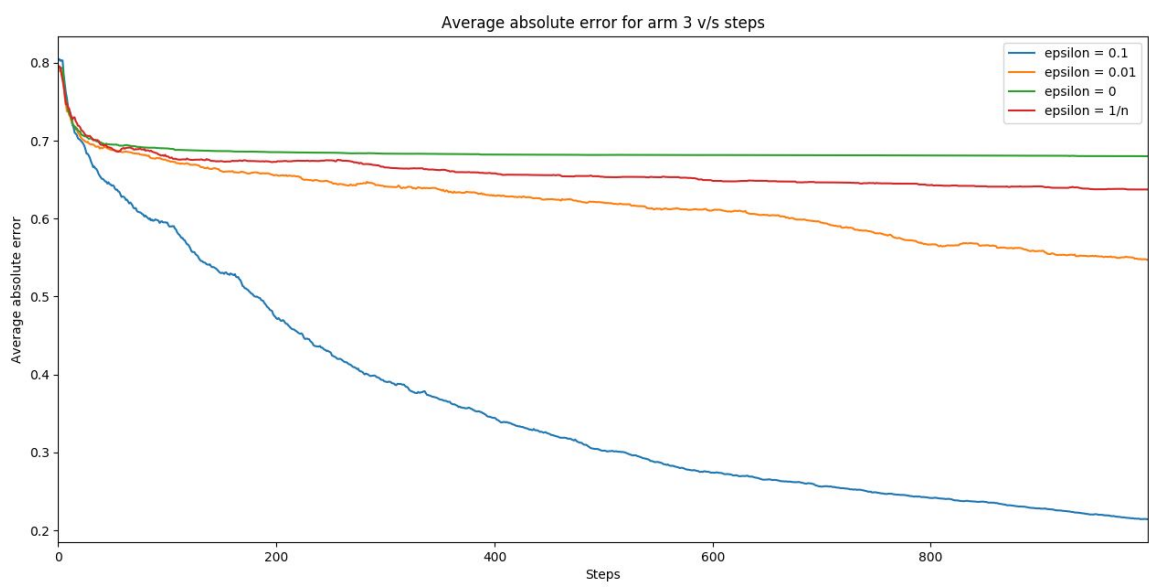
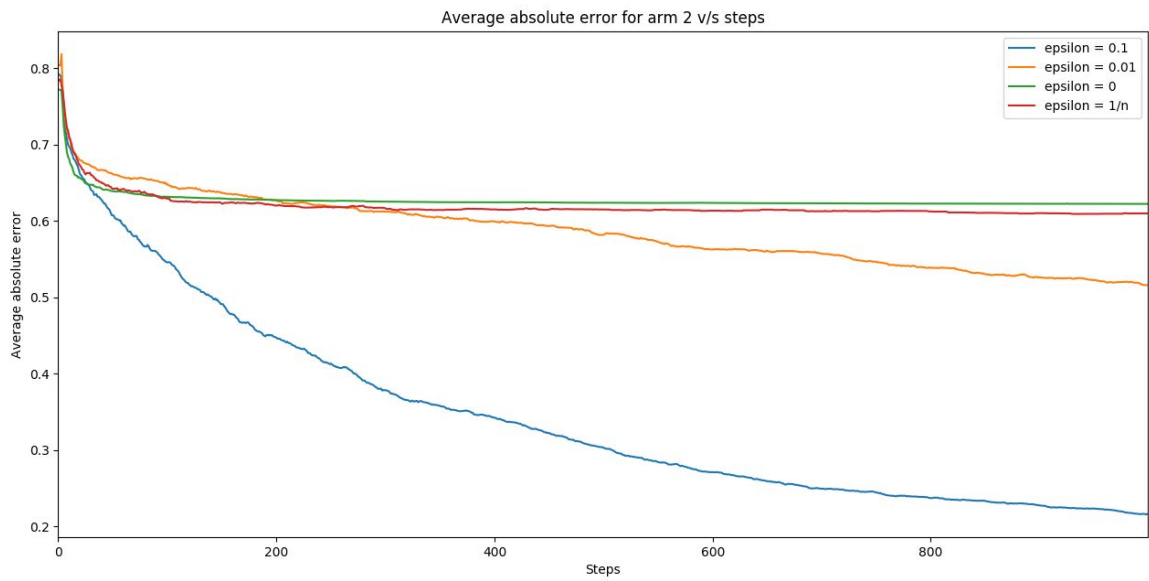
Mean of 10 arms were picked from normal distribution with mean = 0, variance = 1. Variance of each arm was 1. Training was done using epsilon-greedy. Experiment was averaged over 2000 iterations each running for 1000 steps.

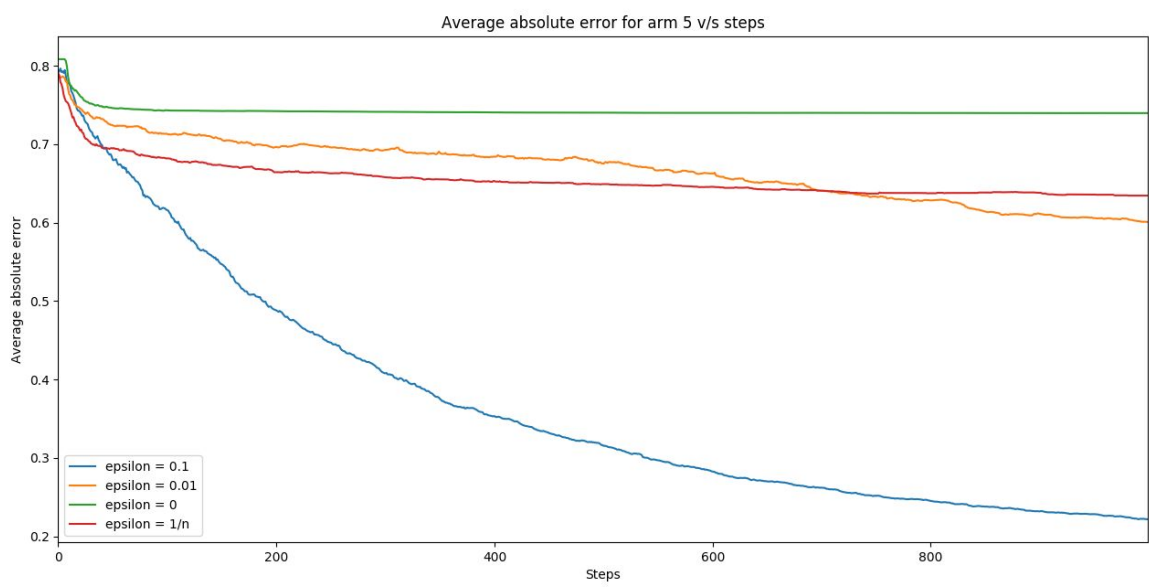
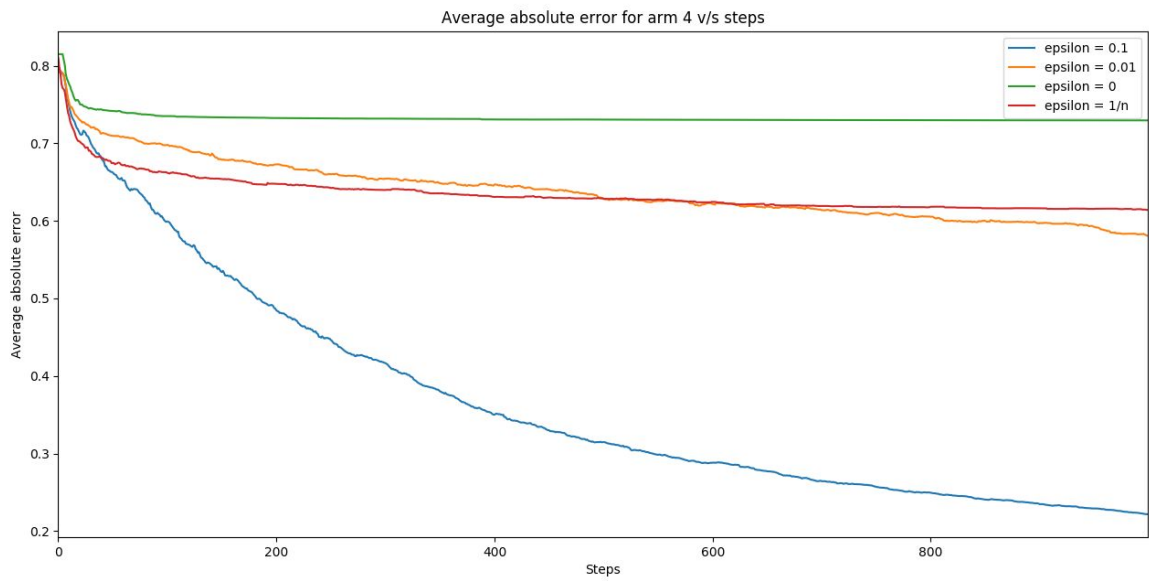


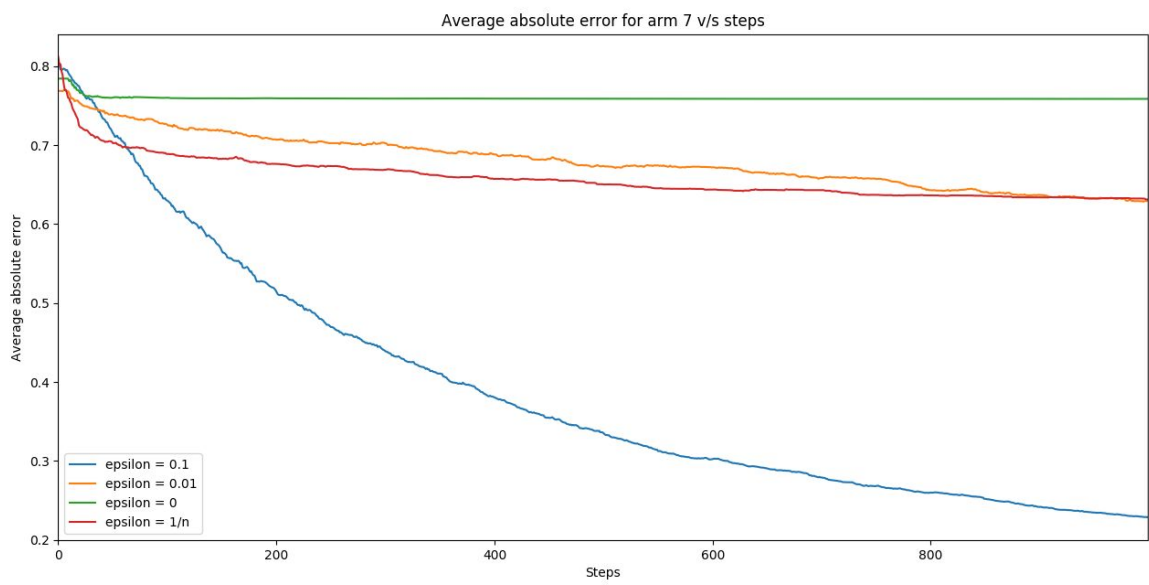
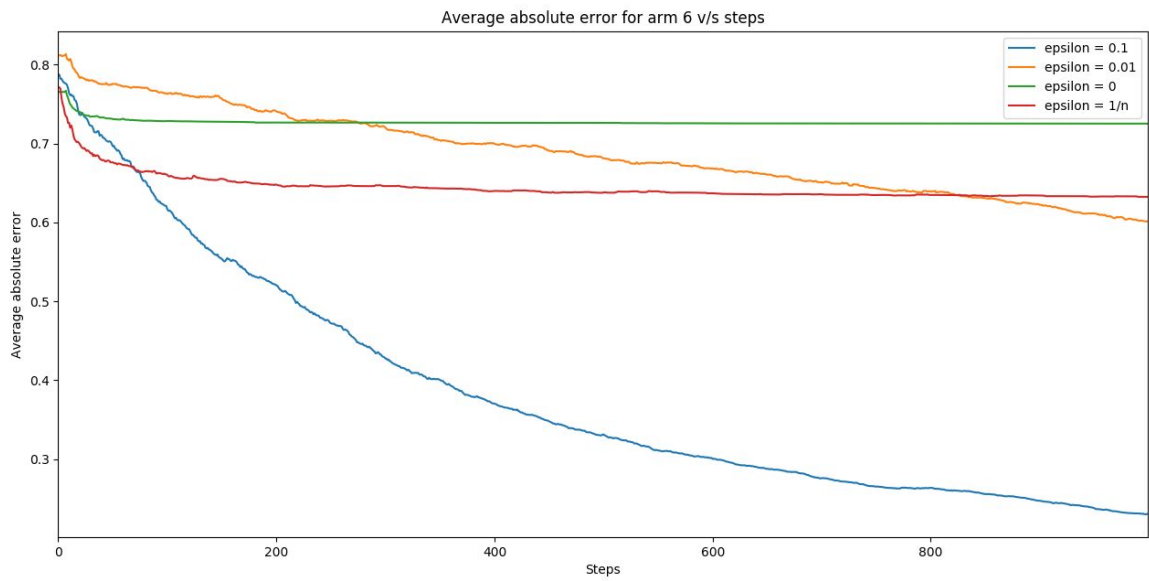


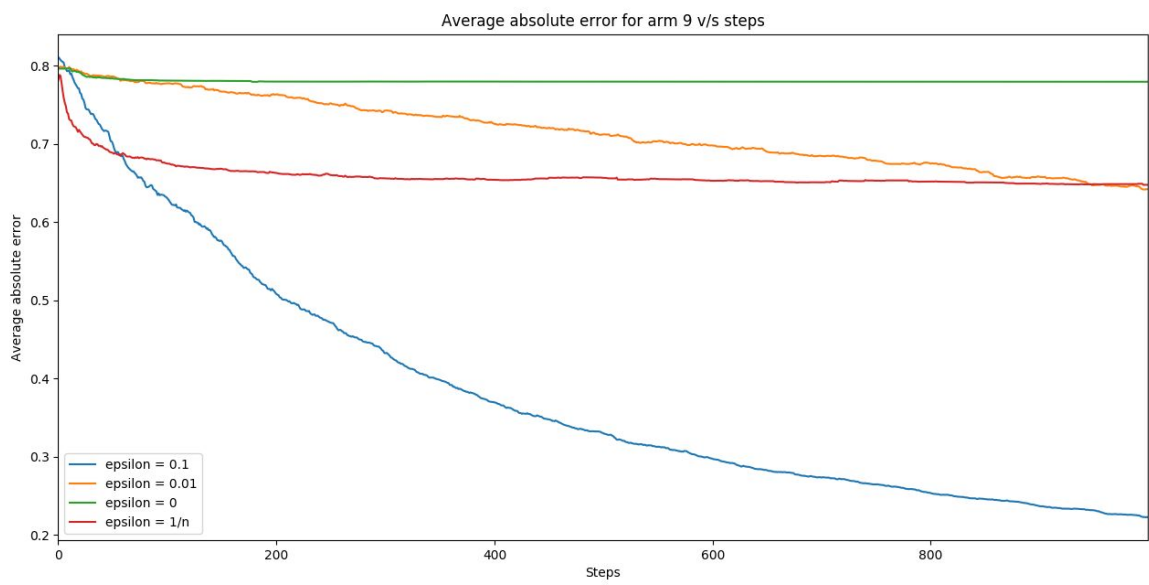
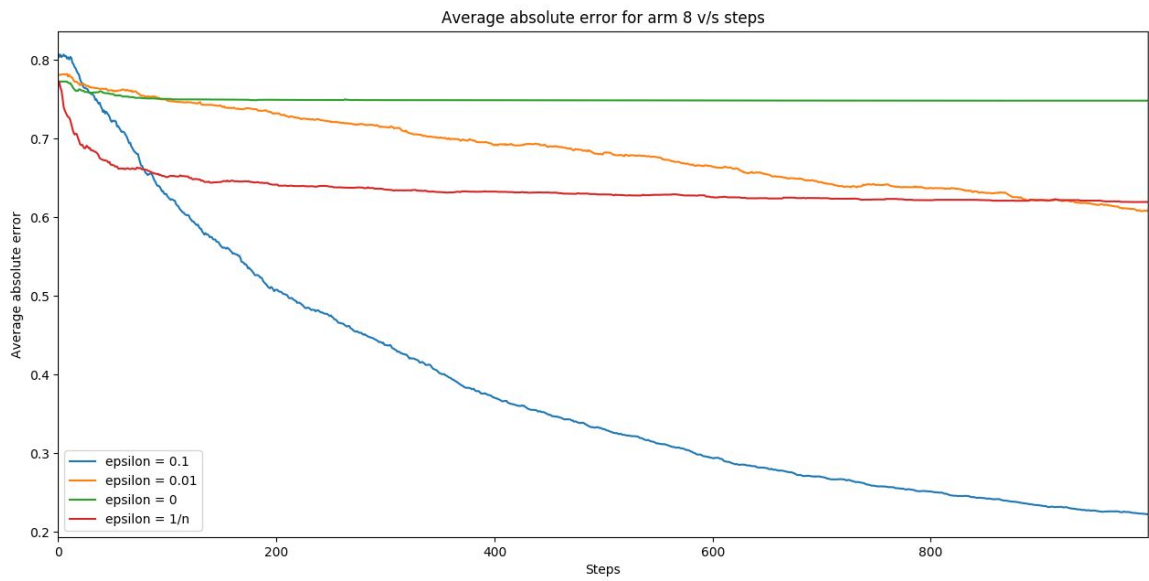
The following plots are for each arm numbered from 0 to 9 inclusive. They show the average absolute error of each arm across 2000 experiments each running for 1000 steps (same as before).





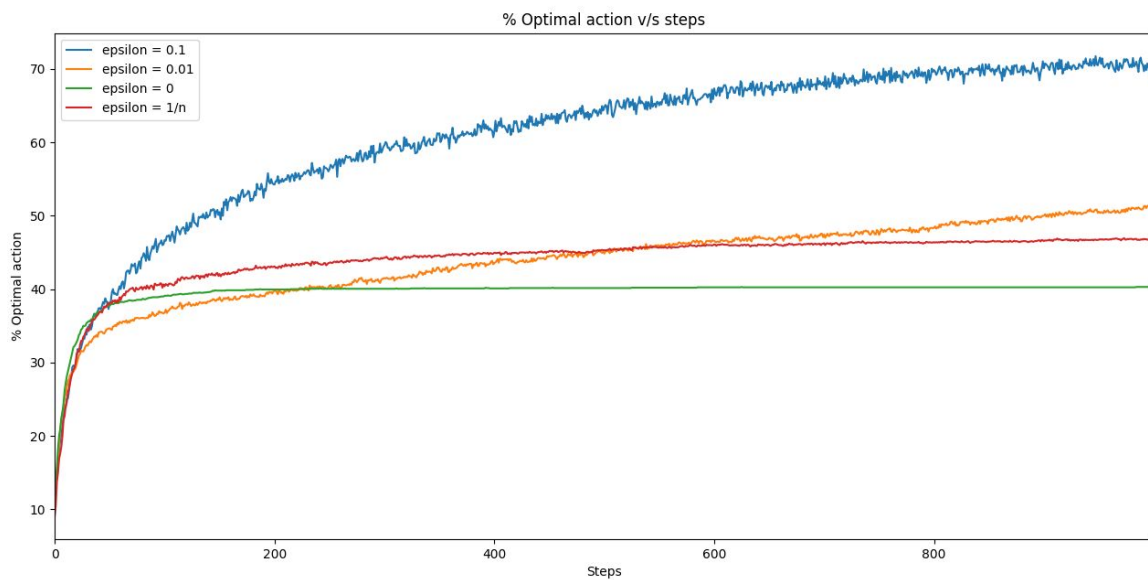
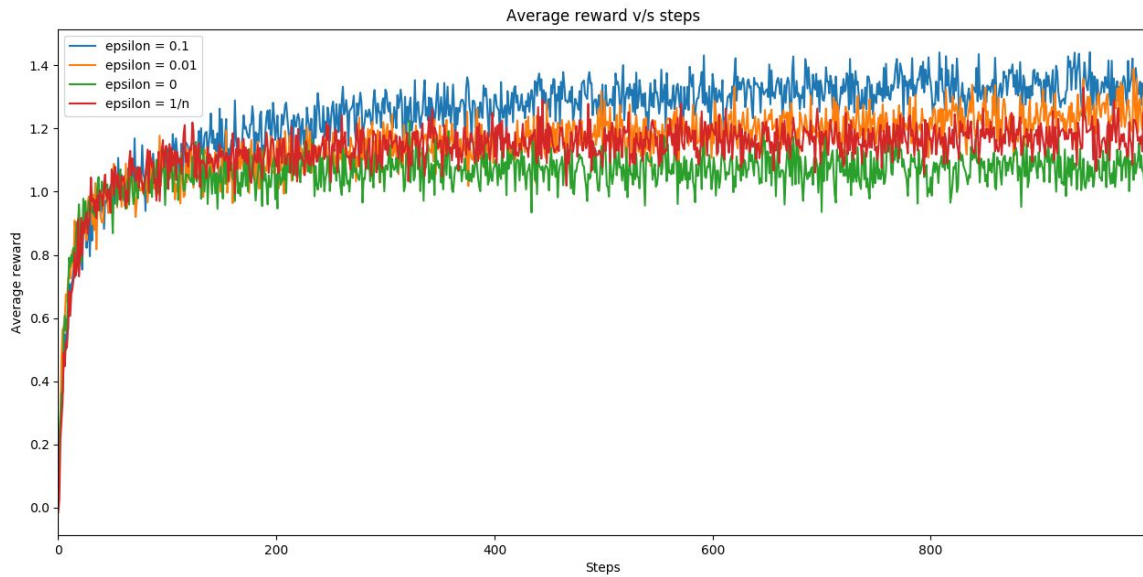




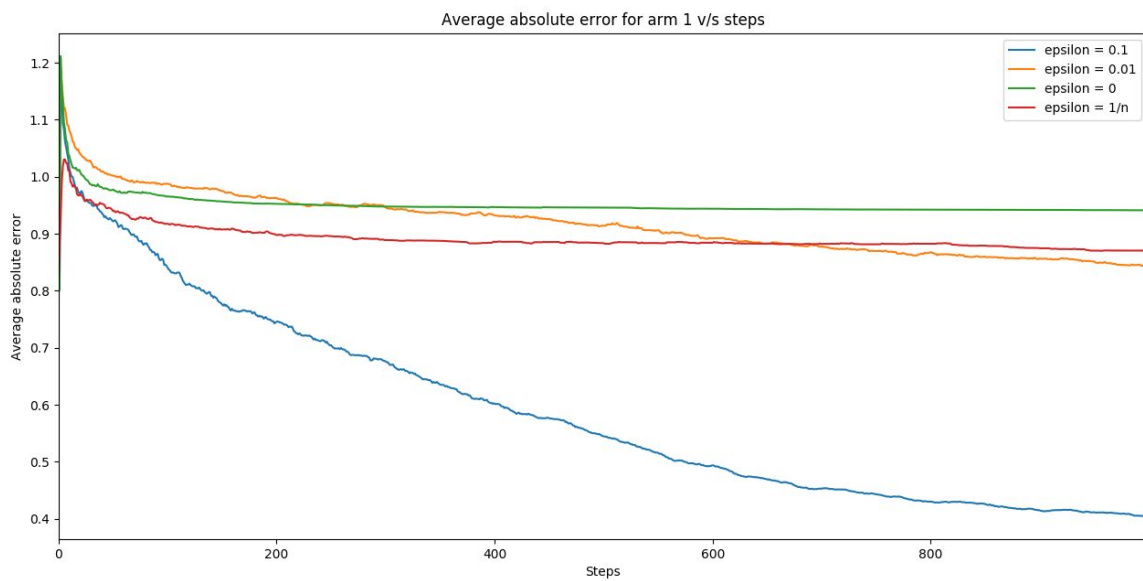
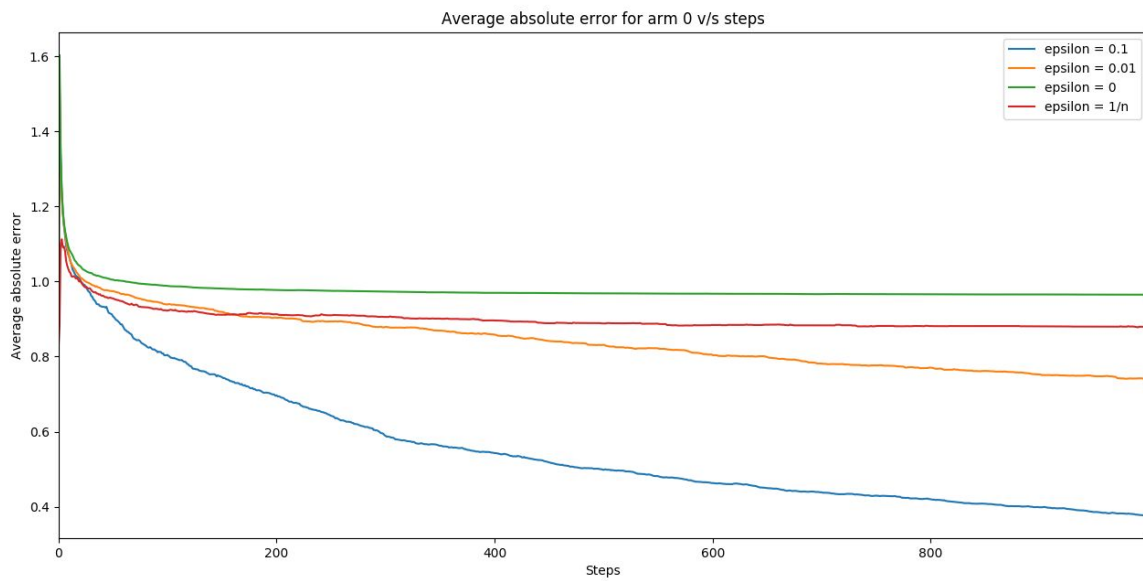


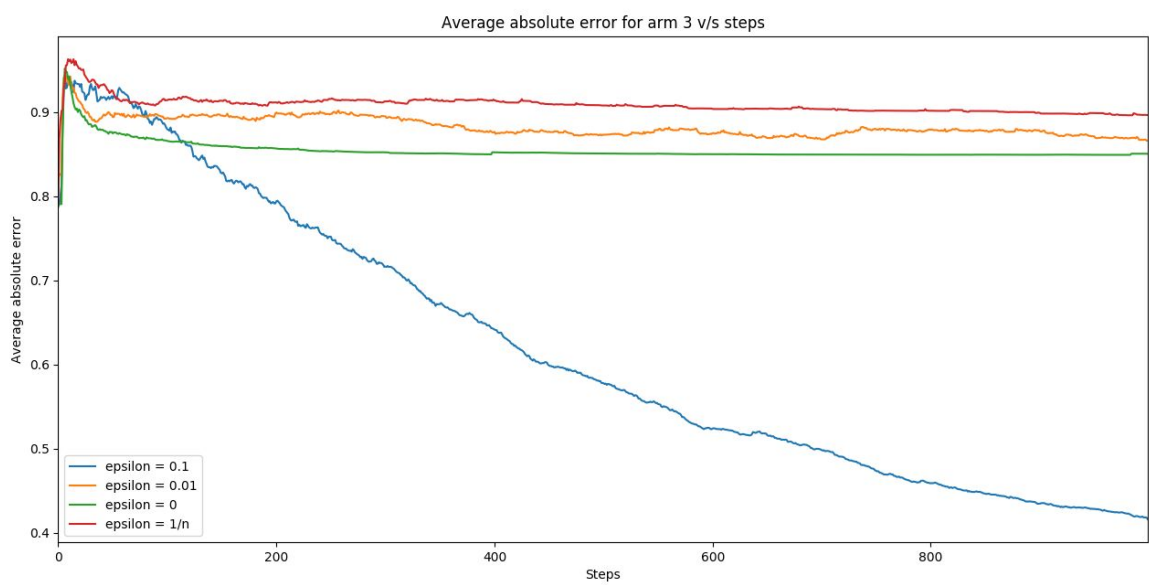
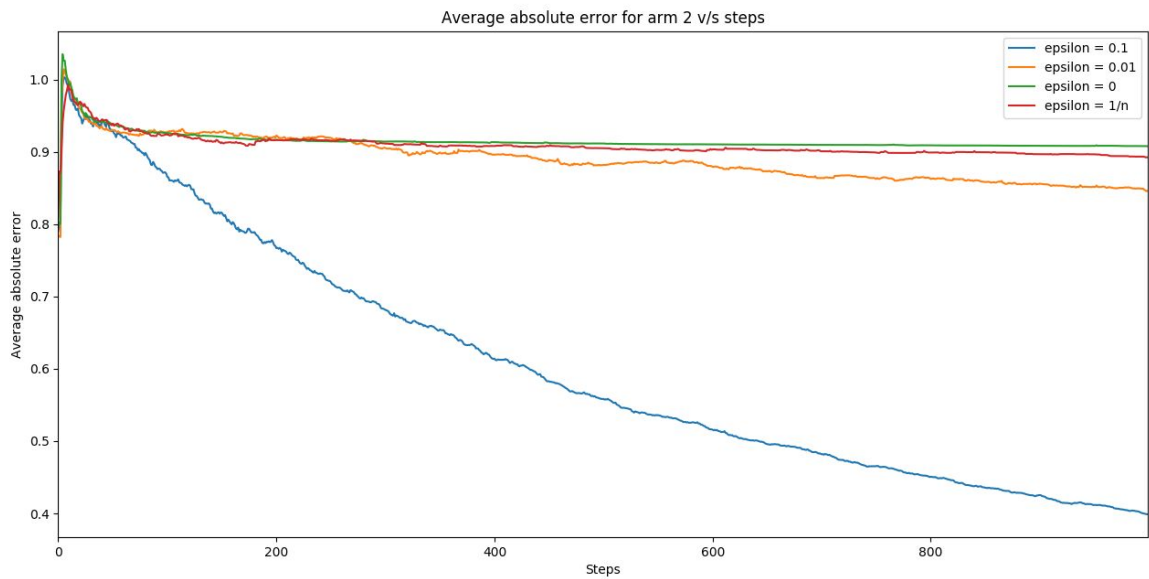
Q2.) $\epsilon = 1/n$ is a sequence of epsilon which satisfies Eq 2.7 of SB. Here n is the number of current timestep.

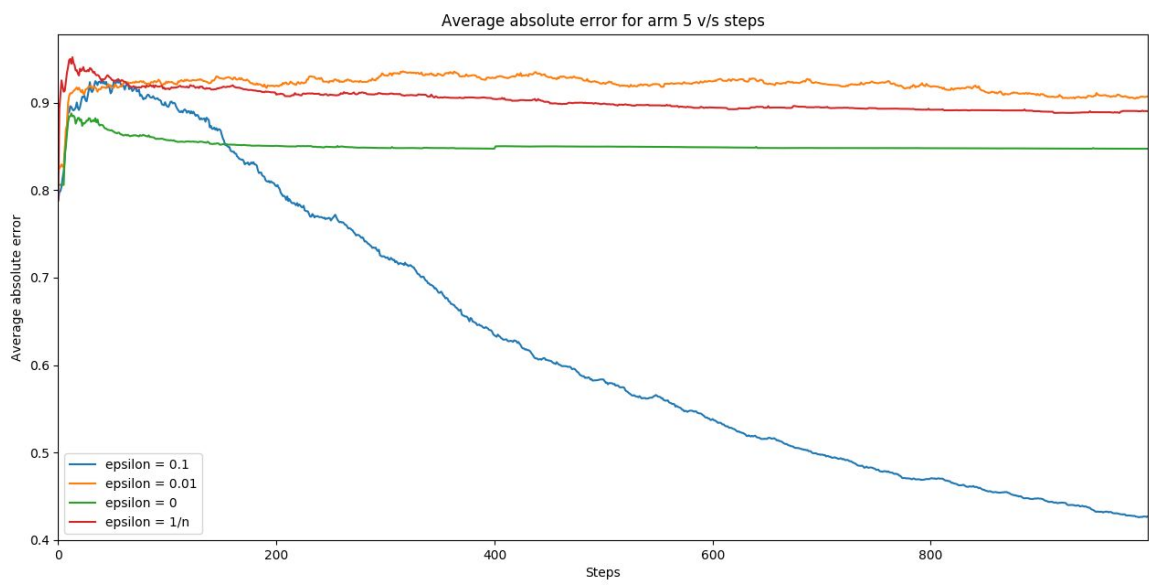
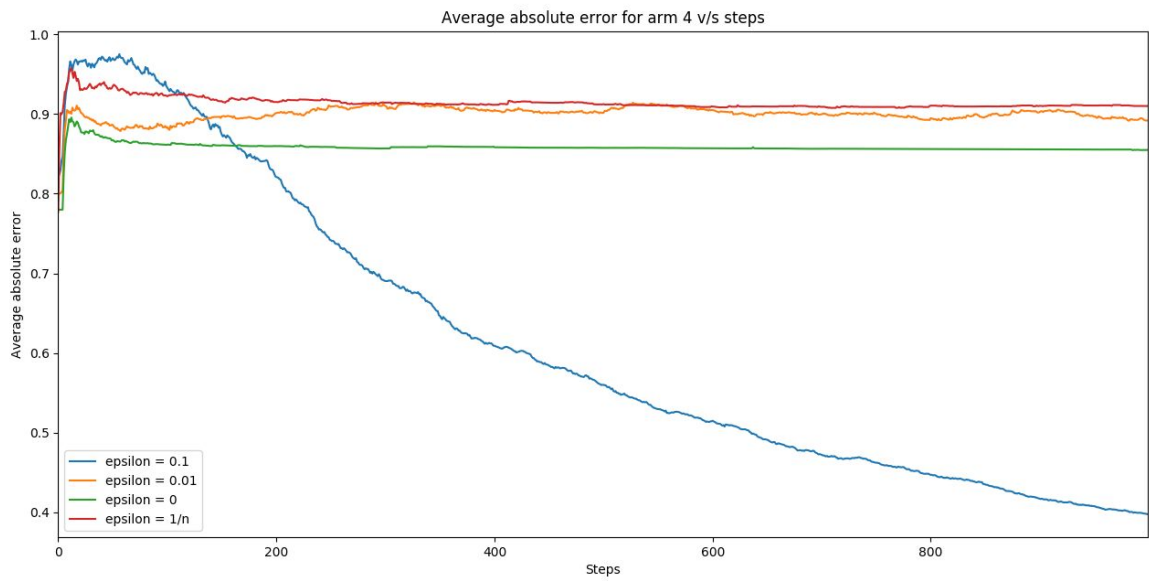
Mean of 10 arms were picked from normal distribution with mean = 0, variance = 1. Variance of each arm was 4. Training was done using epsilon-greedy. Experiment was averaged over 2000 iterations each running for 1000 steps.

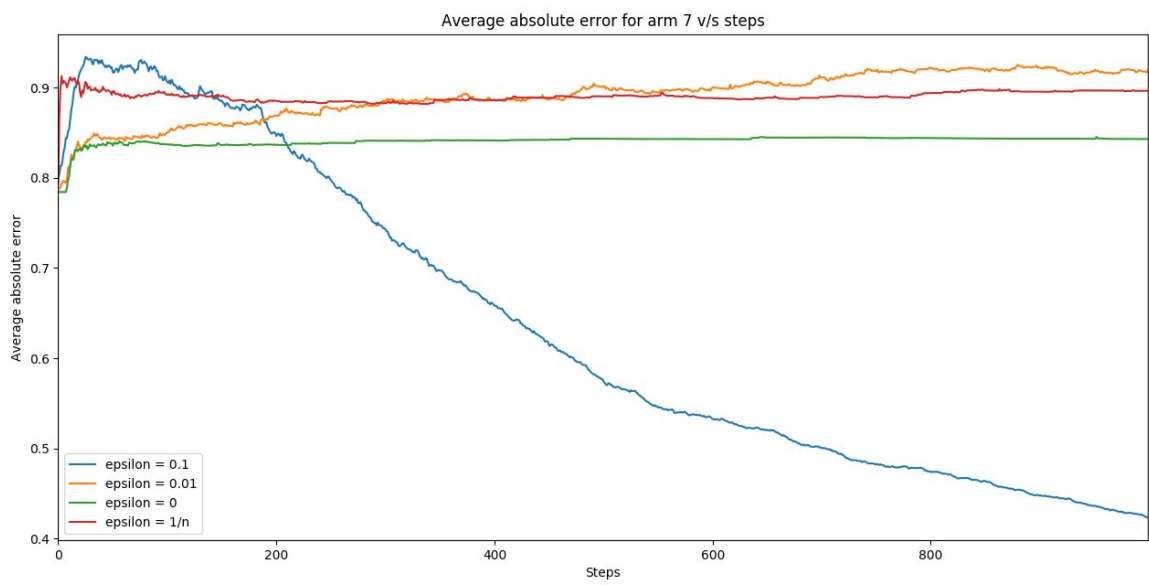
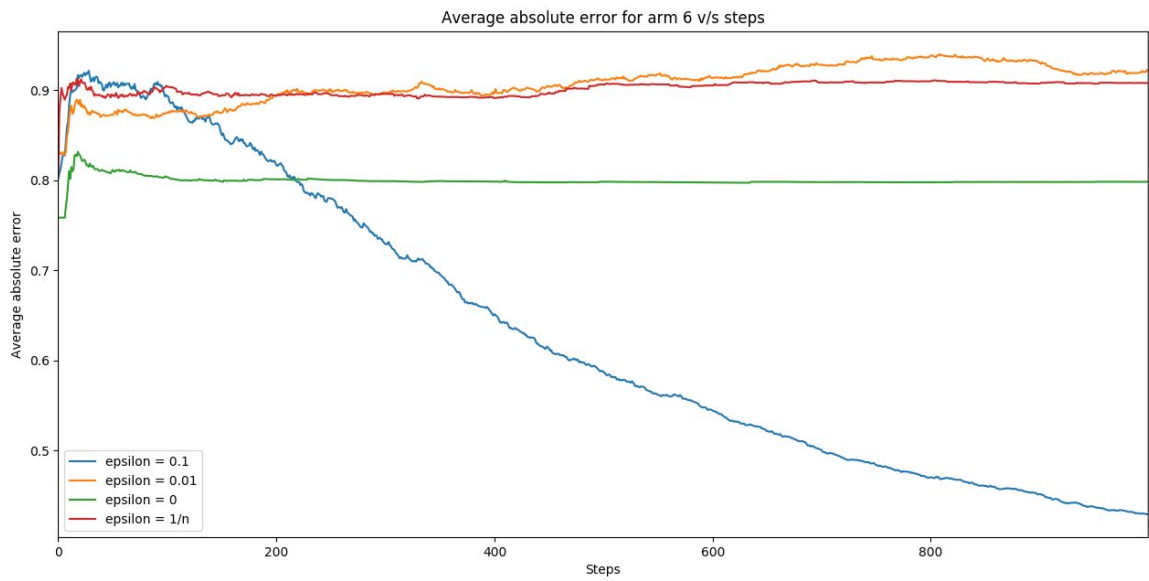


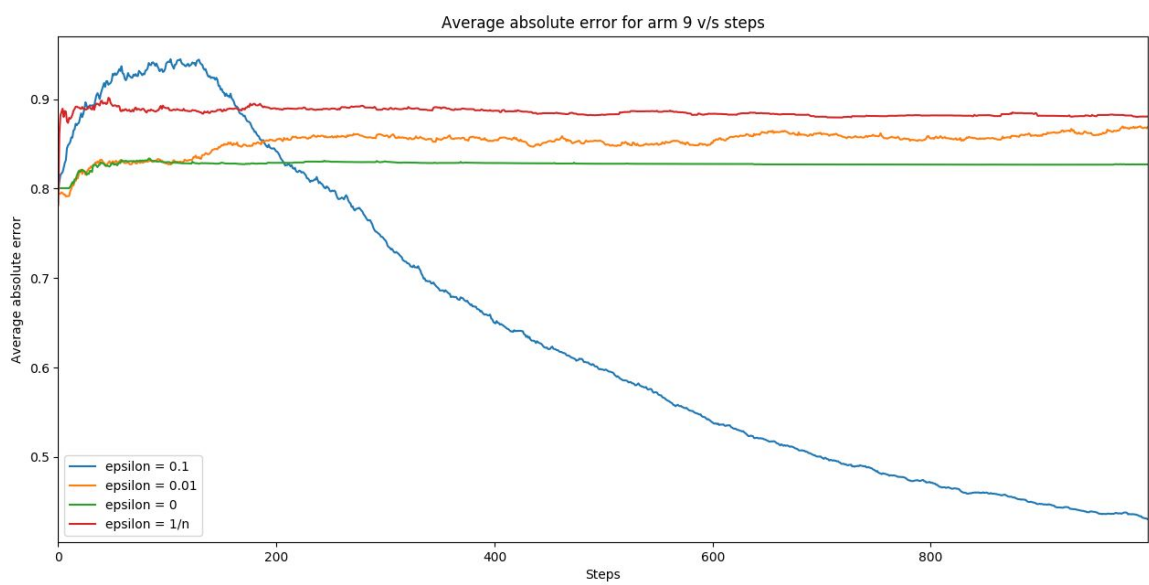
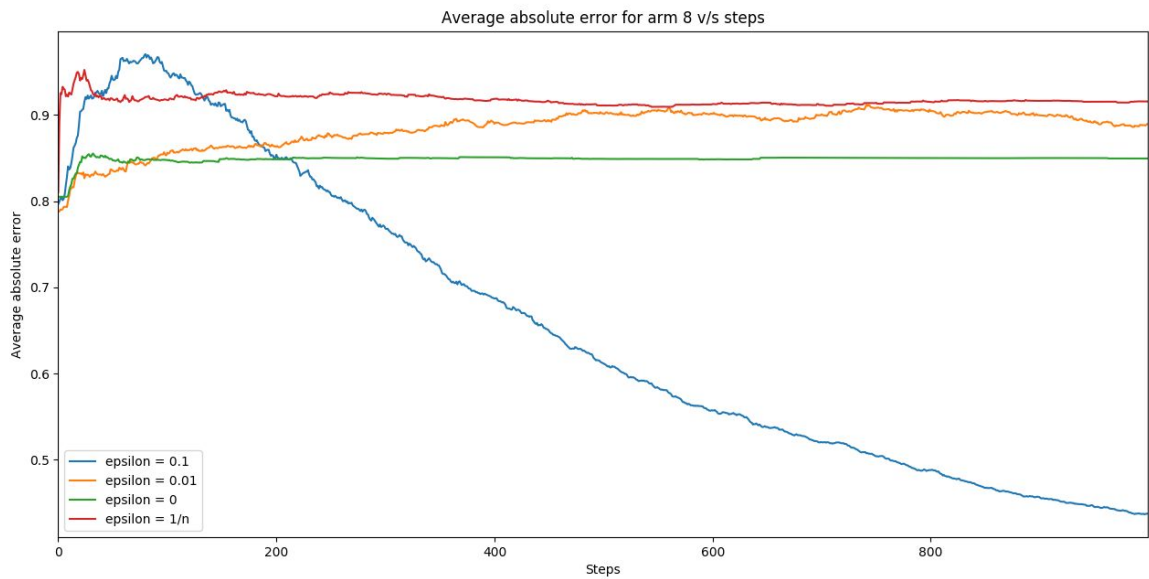
The following plots are for each arm numbered from 0 to 9 inclusive. They show the average absolute error of each arm across 2000 experiments each running for 1000 steps (same as before).



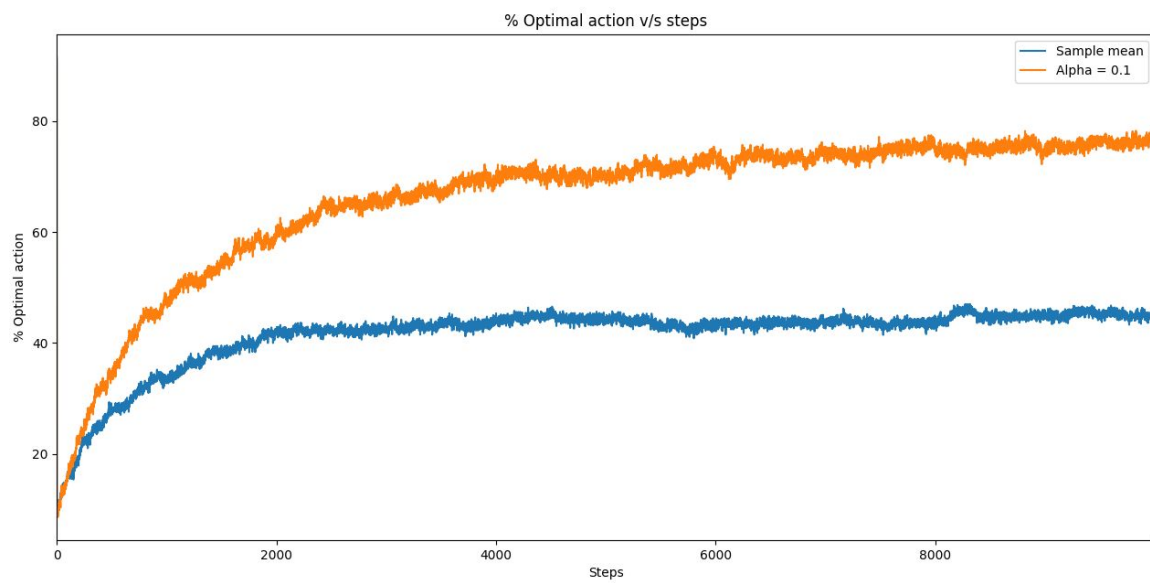
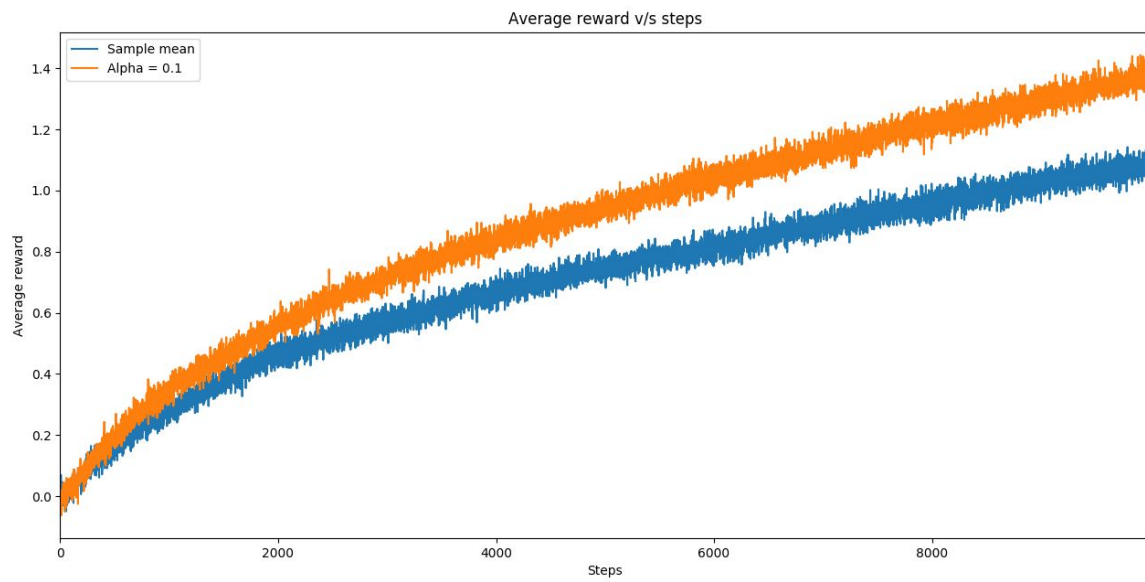




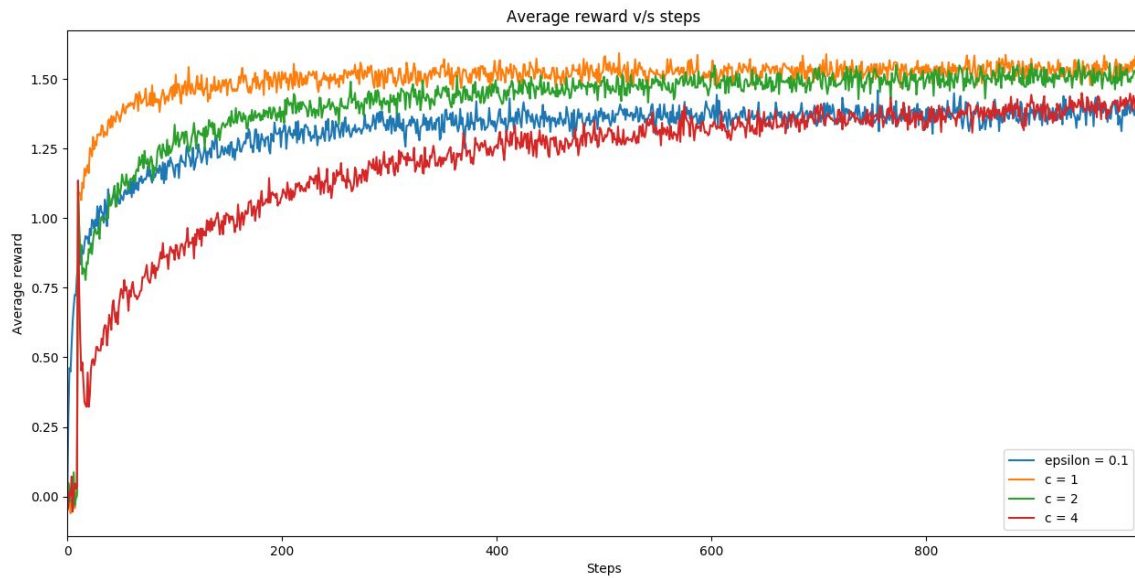




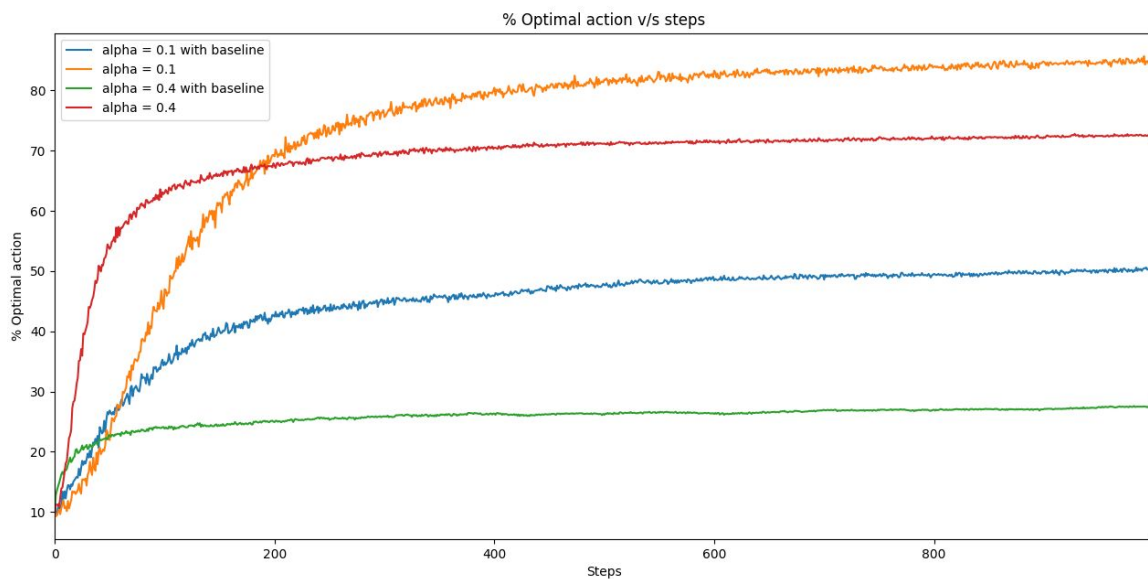
Q5.) The following plots are for non stationary problems as described in Ex 2.5 of SB



Q6.) Average performance of UCB action selection on the 10-armed testbed. The results are for $c = 1, 2, 4$ and $\epsilon = 0.1$ for epsilon greedy.



Q7.)



Theory questions follow...

7/9/20

Sem - 5

Monsoon 2020

Reinforcement Learning

HW 1

Setu Gupta (2018190)

ECE564

Q-3 In the long run the best performing strategy is the one which has

$$P(E) = P(\text{Choosing the optimal arm as } \text{step} \rightarrow \infty)$$

highest. This strategy will give highest cumulative reward in long run.

a) $\epsilon = 0$ (Purely greedy)

This can possibly give $P(E)$ but there is a high possibility that we get stuck at some sub optimal choice of arm.

b) $\epsilon = \text{Constant}$

$$P(E) = 1 - \epsilon + \frac{\epsilon}{|A|} = 1 + \epsilon \left(\frac{1}{|A|} - 1 \right)$$

Here $\frac{1}{|A|} - 1$ is Constant $= \frac{1}{10} - 1 = 0.1 - 1 = -0.9$

$$\therefore P(E) = 1 - 0.9\epsilon$$

\therefore lower ϵ performs better $\therefore 0.01$

$$\begin{aligned} P(E | \epsilon = 0.01) &= 1 - (0.01 \times 0.9) \\ &= 1 - 0.009 \\ &= 0.991 \end{aligned}$$

$$c.) \quad \epsilon = 1/n$$

$$P(\epsilon) = \lim_{n \rightarrow \infty} 1 - 0.9 \left(\frac{1}{n} \right) = 1$$

\therefore This strategy always picks the best arm as $n \rightarrow \infty$.

\therefore The best performing strategy among 4 options is $\epsilon = 1/n$.

$$Q-4) \quad Q_{n+1} = Q_n + \beta (R_n - Q_n)$$

where β is step size

a) Sample mean i.e. $\beta = 1/n$.

$$\therefore Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

+ a (omitting writing $Q_{n+1}(a)$ for clarity)

Consider Q_2 i.e. $n=1$

$$\begin{aligned} Q_2 &= Q_1 + \frac{1}{1} [R_1 - Q_1] \\ &= \cancel{Q_1} + R_1 - \cancel{Q_1} \\ &= R_1 \end{aligned}$$

\therefore After the action is chosen at least once, the dependence on Q_1 goes away.

b) (i) Constant step-size i.e. $\beta = \alpha$

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

$$= \alpha_n R_n + (1-\alpha) Q_n$$

$$= \alpha_n R_n + (1-\alpha) [\alpha_{n-1} R_{n-1} + (1-\alpha) Q_{n-1}]$$

$$= \alpha_n R_n + (1-\alpha) \alpha_{n-1} R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \alpha_n R_n + (1-\alpha) \alpha_{n-1} R_{n-1} + (1-\alpha)^2 \alpha_{n-2} R_{n-2} \dots$$

$$\dots + (1-\alpha)^{n-1} \alpha_1 R_1 + (1-\alpha)^n Q_1$$

$$= \alpha_n \sum_{i=0}^{n-1} (1-\alpha)^i \alpha_{n-i} R_{n-i} + (1-\alpha)^n Q_1$$

\therefore Dependent on Q_1 ,

(ii) Notice that $0 < \alpha \leq 1$

\therefore for smaller α , $(1-\alpha)^n$ is larger i.e. larger dependence on Q_1 .

(iii) I propose two way of dealing w/ dependence on Q_1 ,

MI $\alpha = 1$

$$\begin{aligned} \Rightarrow Q_{n+1} &= Q_n + 1 \cdot [R_n - Q_n] \\ &= \cancel{Q_n} + R_n - \cancel{Q_n} = R_n \end{aligned}$$

Drawback: No regard to history

MI Instead of choosing values of $Q_1(a)$ $\forall a$ manually, pull each arm one by one to obtain

$R_0(a) \forall a$

Set $Q_1(a) = R_0(a) \forall a$

This way we don't need to decide the values.

Drawback: Does not follow constant α strategy for first $|A|$ actions.

$$Q-6) A_t = \operatorname{argmax}_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Note that initially $Q_t(a) = \lambda$ (constant) $\forall a$

Also $N_t(a) = 0 \forall a$

\therefore For the first step the value

$$Q(a) \triangleq \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right] \rightarrow \infty$$

\therefore First action is chosen at random. Suppose we

Choose action i st. $1 \leq i \leq |A|$

$\Rightarrow Q_2(i)$ becomes finite as $N_2(i) = 1$

\therefore We will not choose i this time.

By the same logic we can say that for first 10 steps we choose every action once in an arbitrary order.

Now let's look @ 11th step.

$$c \sqrt{\frac{\ln t}{N_t(a)}} = c \sqrt{\frac{\ln(11)}{1}} \quad \forall a$$

\therefore This choice will purely be influenced by

$Q_{t+1}(a) \forall a.$

Now across 2000 runs, The probability of $Q_{11}(a^*)$ being higher for truly optimal action a^* is higher

\therefore (a) 11th step since $P[Q_{11}(a^*) > Q_{11}(a)]$
 $\forall a \in A, a \neq a^*$

is high, we pick a^* which explains the spike.

(b) 12th step $S_{12}(a^*) = \left[Q_{12}(a^*) + c \sqrt{\frac{\ln(12)}{2}} \right]$

Note that the factor $c \sqrt{\frac{\ln(12)}{N_t(a)}}$ $\forall a$ is lowest

for a^* as $N_t(a^*) = 2$ whereas $N_t(a) \forall a \neq a^*$ is 1, \therefore it is $1/\sqrt{2}$ times less for a^* .

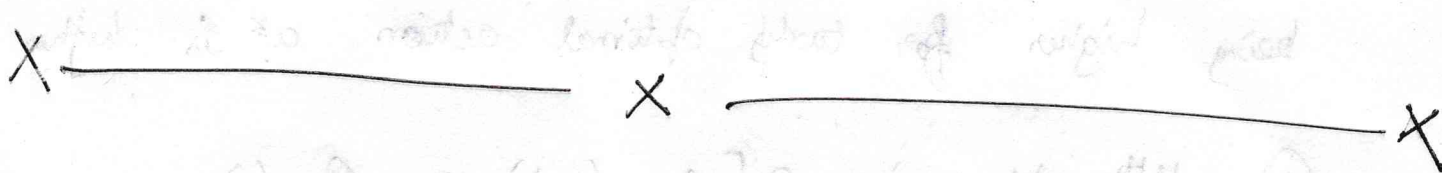
Now this low value may be overpowered by $Q_{12}(a^*)$

but if c is less, this is ~~more~~ unlikely. Hence we get a dip.

This is also the reason that the spike is more pronounced for higher c as higher the c , more will $c \sqrt{\frac{\ln(t)}{N(t)}}$ overpower $Q_t(a)$.

Also note that it is most pronounced (a) 11th step as

at later steps, values of $S(a)$ & $T(a)$ start to converge.



$$[x_0 + a, x_0 + a + \dots] \quad (1)$$

at later steps, values of $S(a)$ & $T(a)$ start to converge.

$$\left[\frac{(x_0 + a)}{2} + (x_0 + a) \right] = (x_0 + a) \quad (2)$$

$$x_0 + a = \frac{(x_0 + a)}{2} + (x_0 + a) \quad (3)$$

$$x_0 + a = (x_0 + a) \quad (4)$$

$$x_0 + a = (x_0 + a) \quad (5)$$

$$(x_0 + a) = (x_0 + a) \quad (6)$$

$$(x_0 + a) = (x_0 + a) \quad (7)$$

$$(x_0 + a) = (x_0 + a) \quad (8)$$

$$(x_0 + a) = (x_0 + a) \quad (9)$$

$$(x_0 + a) = (x_0 + a) \quad (10)$$

$$(x_0 + a) = (x_0 + a) \quad (11)$$

$$(x_0 + a) = (x_0 + a) \quad (12)$$