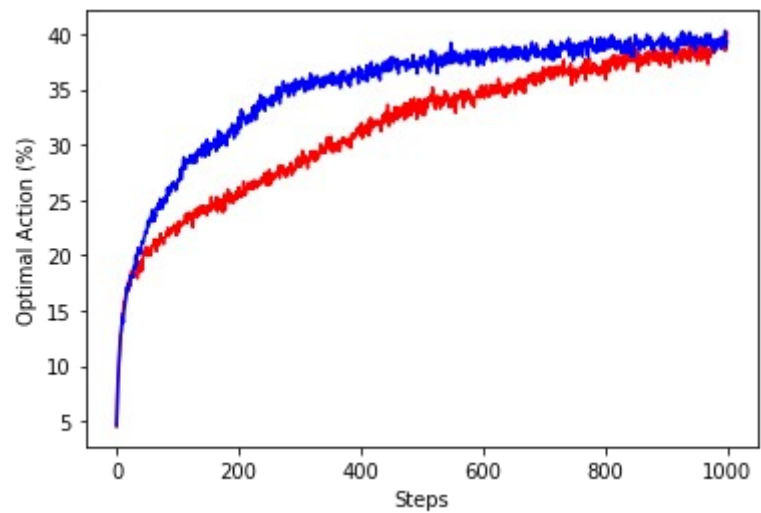
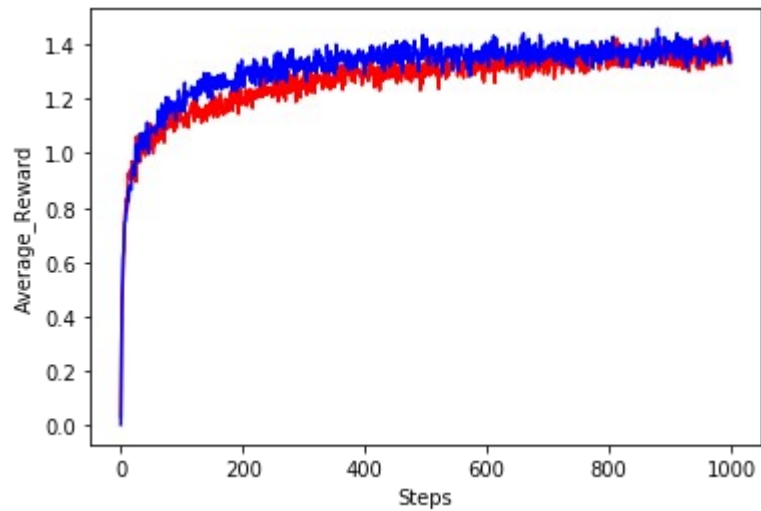


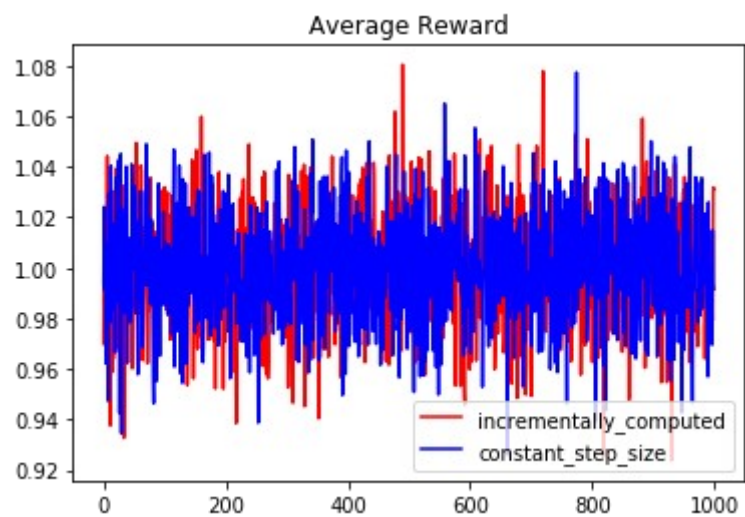
Assignment 1

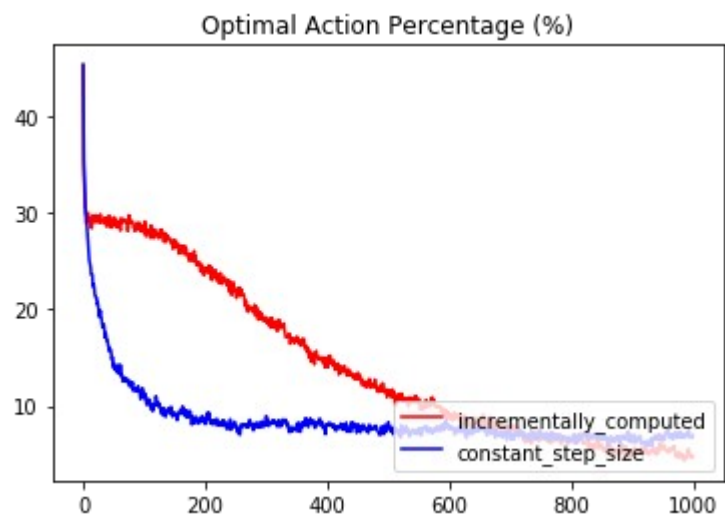
Question 1

`q = np.random.normal(0,1,n_arms)`



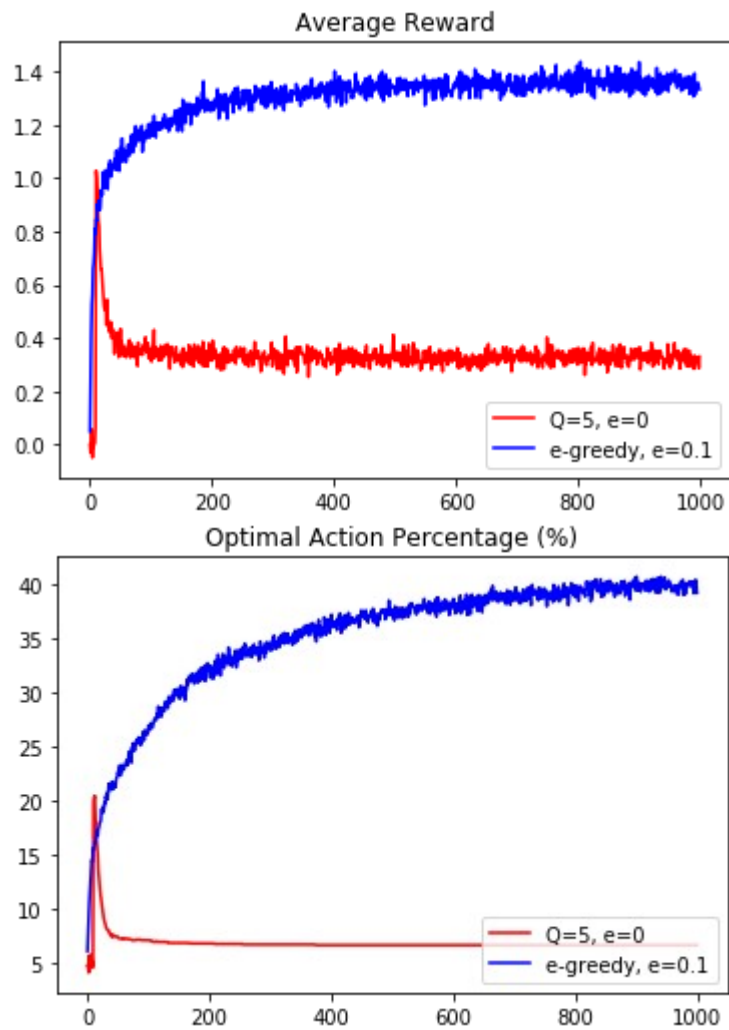
`q = np.ones(n_arms)`



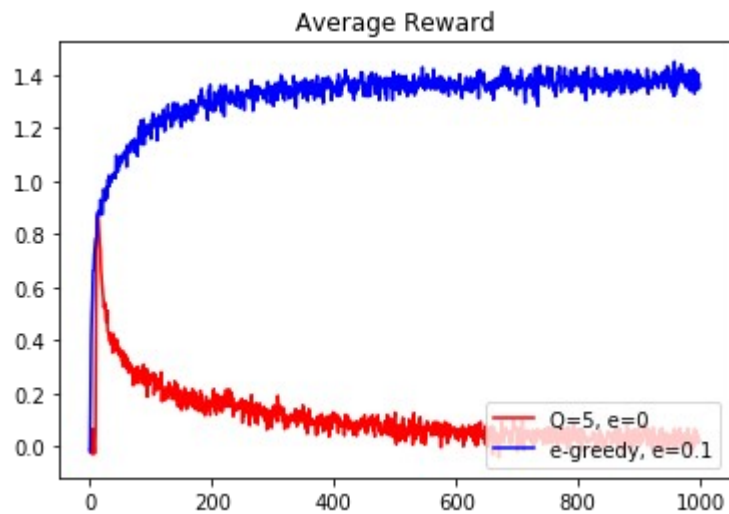


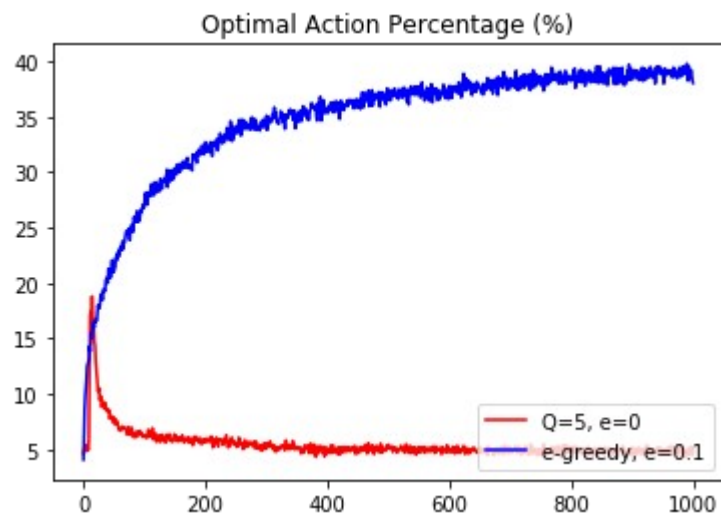
Question 2

Question 2 fig 2.3 for stationary ; incrementally computed

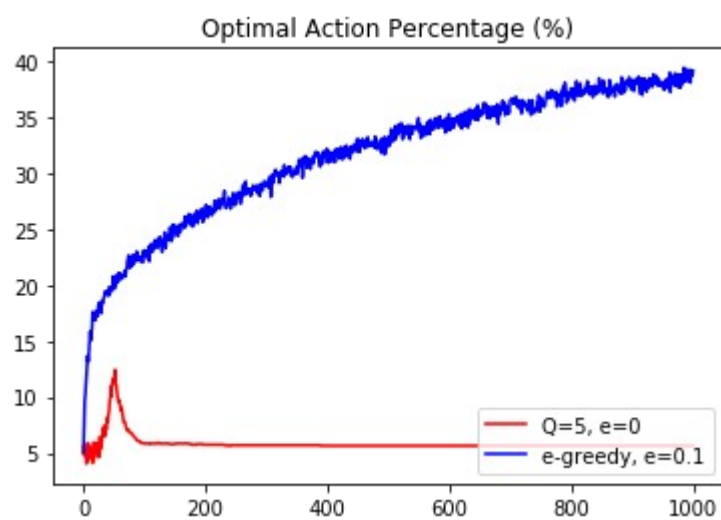
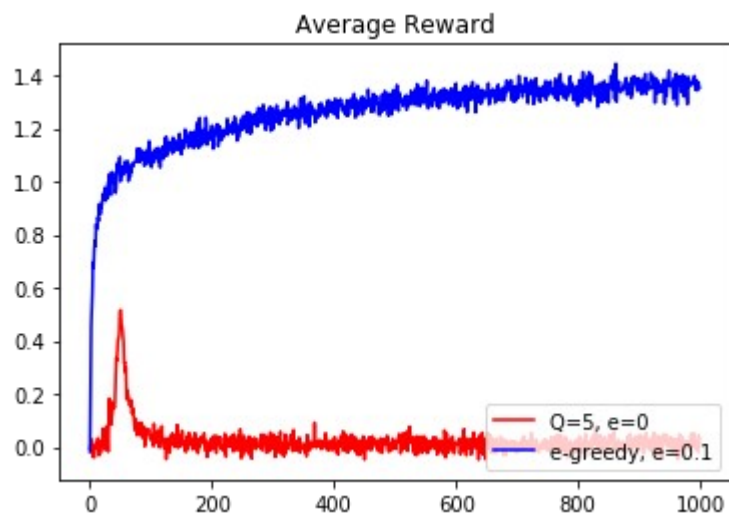


Question 2 fig 2.3 for ; non-stationary; incrementally-computed

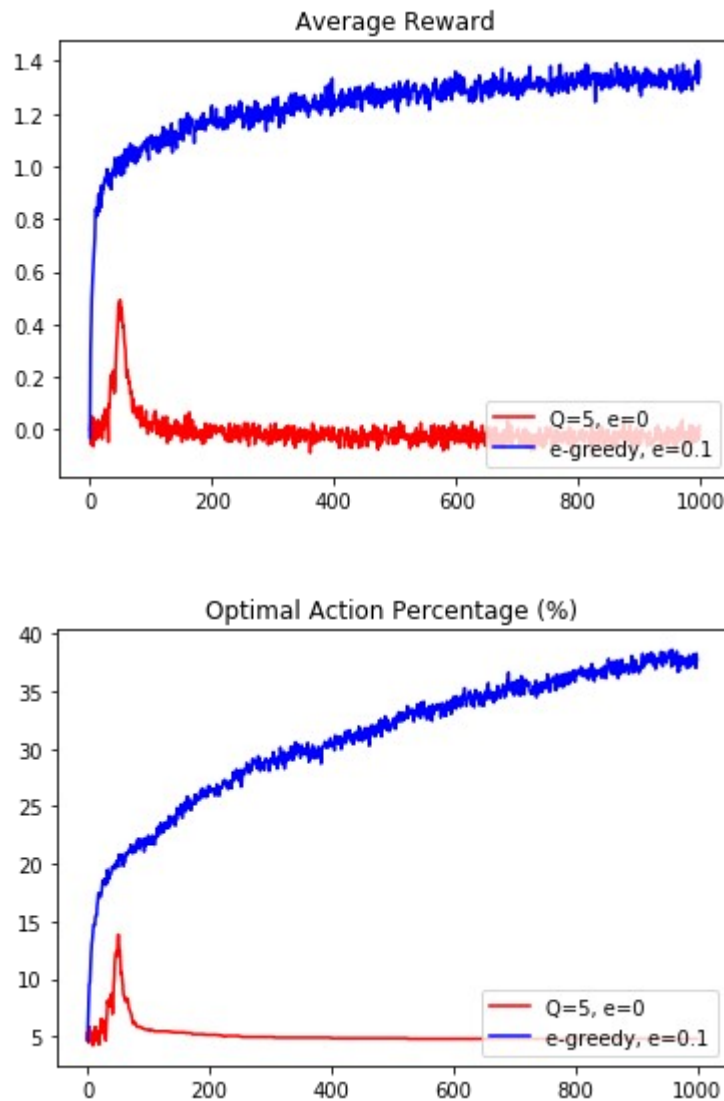




Question 2 fig 2.3 for stationary ; constant_step-size



Question 2 fig 2.3 for non stationary ; constant_step-size



Initially, the optimistic method performs worse because it explores more, but eventually it performs better because its exploration decreases with time.

Solution for question 6

The oscillations in start of optimistic curve show that estimated values are still optimistic. When the curve increases, it is due to choosing the bandit with the highest expected reward. Estimated values are still optimistic (much higher than their expected values). In this context, the best bandit is selected due to its optimistic value during the early part of the curve, but when its estimated value becomes much closer to the real expected value, other decisions with optimistic values are chosen.

Due to this the curve oscillates during the early stages, until all estimated values aren't as nearly as optimal and become much closer to the real expected value instead.

Thus optimistic greedy is not better for small number of time steps. Realistic e-greedy performs better for that. Optimistic method encourage exploration only during the beginning of execution, and would just be useful for stationary problems.

Question 3:

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

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

$$= \frac{\alpha}{\sigma_n} R_n + \left(1 - \frac{\alpha}{\sigma_n}\right) Q_n$$

$$= \frac{\alpha}{\sigma_n} R_n + \left(1 - \frac{\alpha}{\sigma_n}\right) \left[\frac{\alpha}{\sigma_{n-1}} R_{n-1} + \left(1 - \frac{\alpha}{\sigma_{n-1}}\right) Q_{n-1} \right]$$

$$= \frac{\alpha}{\sigma_n} R_n + \left(1 - \frac{\alpha}{\sigma_n}\right) \left(\frac{\alpha}{\sigma_{n-1}} \right) R_{n-1} + \left(1 - \frac{\alpha}{\sigma_n}\right)^2 Q_{n-1}$$

$$= \frac{\alpha}{\sigma_n} R_n + \left(1 - \frac{\alpha}{\sigma_n}\right) \left(\frac{\alpha}{\sigma_{n-1}} \right) R_{n-1} + \left(1 - \frac{\alpha}{\sigma_n}\right)^2 \frac{\alpha}{\sigma_{n-2}} R_{n-2} +$$

$$+ \left(1 - \frac{\alpha}{\sigma_n}\right)^{n-1} \frac{\alpha}{\sigma_1} R_1 + \left(1 - \frac{\alpha}{\sigma_n}\right)^n Q_1$$

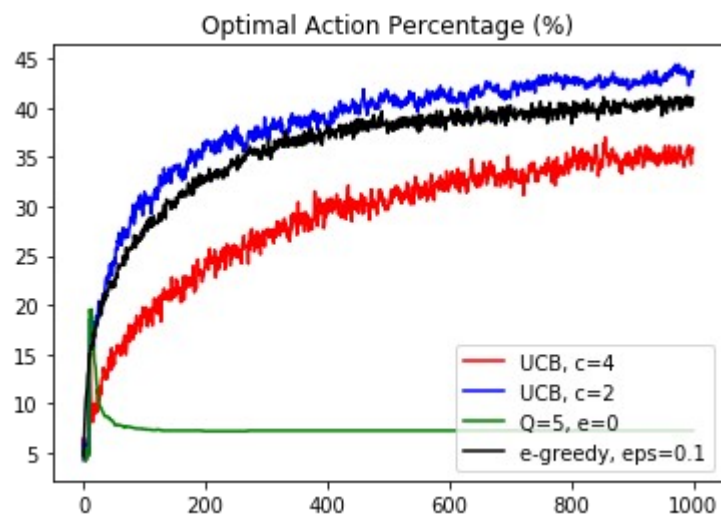
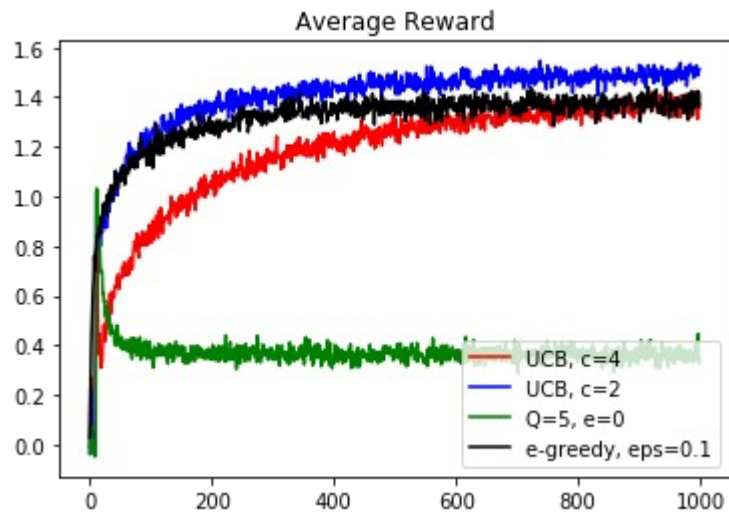
$$= \left(1 - \frac{\alpha}{\sigma_n}\right)^n Q_1 + \sum_{i=1}^n \frac{\alpha}{\sigma_n} \left(1 - \frac{\alpha}{\sigma_n}\right)^{n-i} R_i$$

$$= \frac{1}{\sigma_n} \left[\underbrace{\left(1 - \frac{\alpha}{\sigma_n}\right)^n \sigma_n Q_1}_{=0} + \sum_{i=1}^n \alpha \left(1 - \frac{\alpha}{\sigma_n}\right)^{n-i} R_i \right]$$

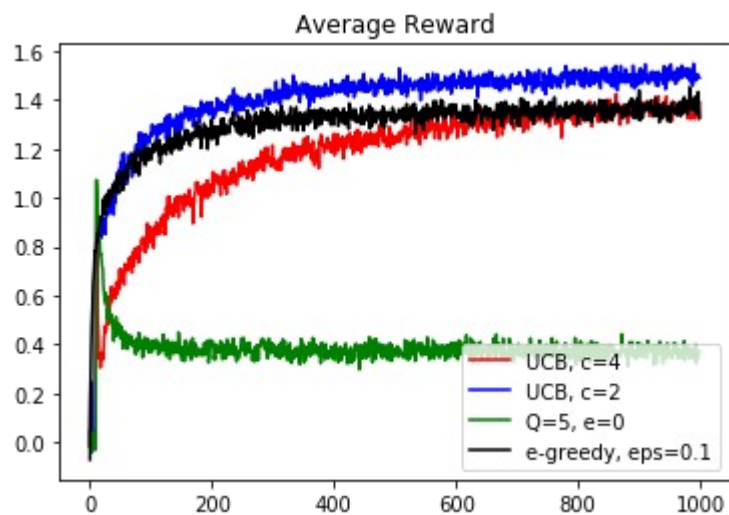
Since $\sigma_0 = 0$ hence $\sigma_n \neq 0$ then $Q_{n+1} = \sum_{i=1}^n \alpha \left(1 - \frac{\alpha}{\sigma_n}\right)^{n-i} R_i$

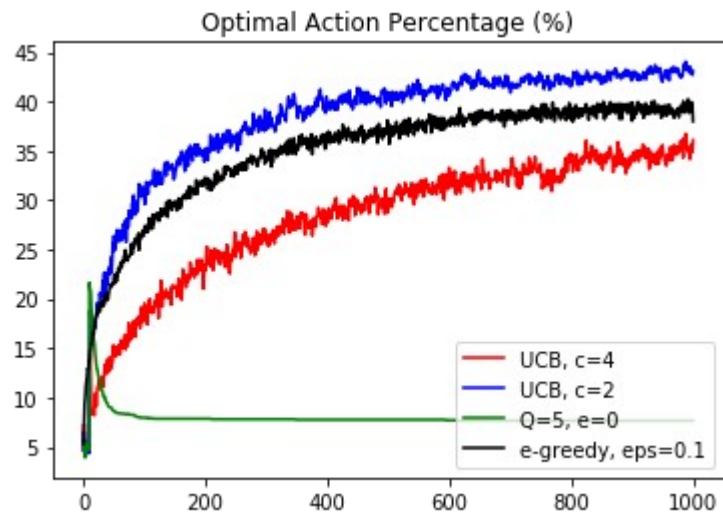
Question 4

UCB for $c=4$, $c=2$, e-greedy = 0.1, $Q=5$, $e=0$; incrementally computed ; stationary

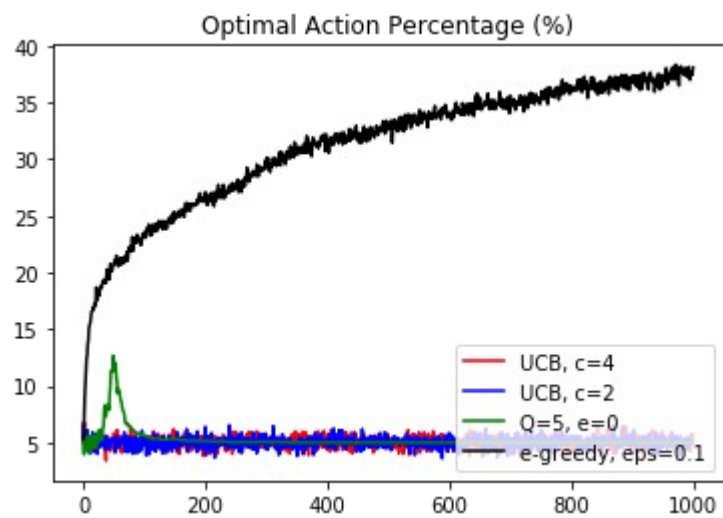
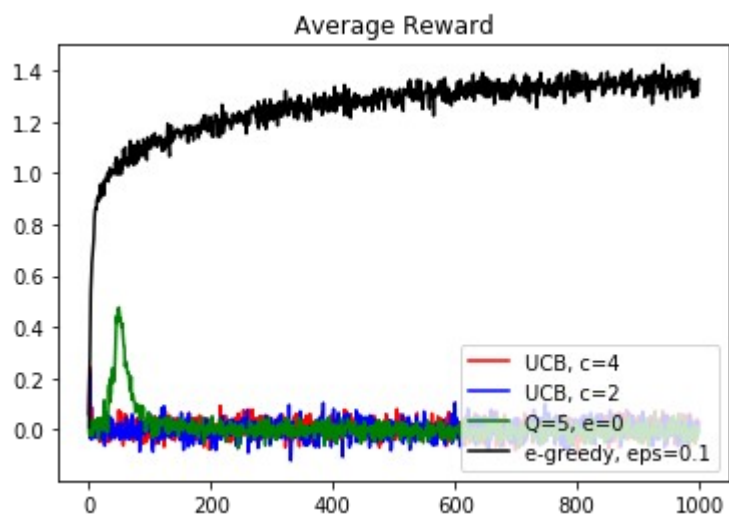


UCB for $c=4$, $c=2$, e-greedy = 0.1, $Q=5$, $e=0$; incrementally computed ; non-stationary

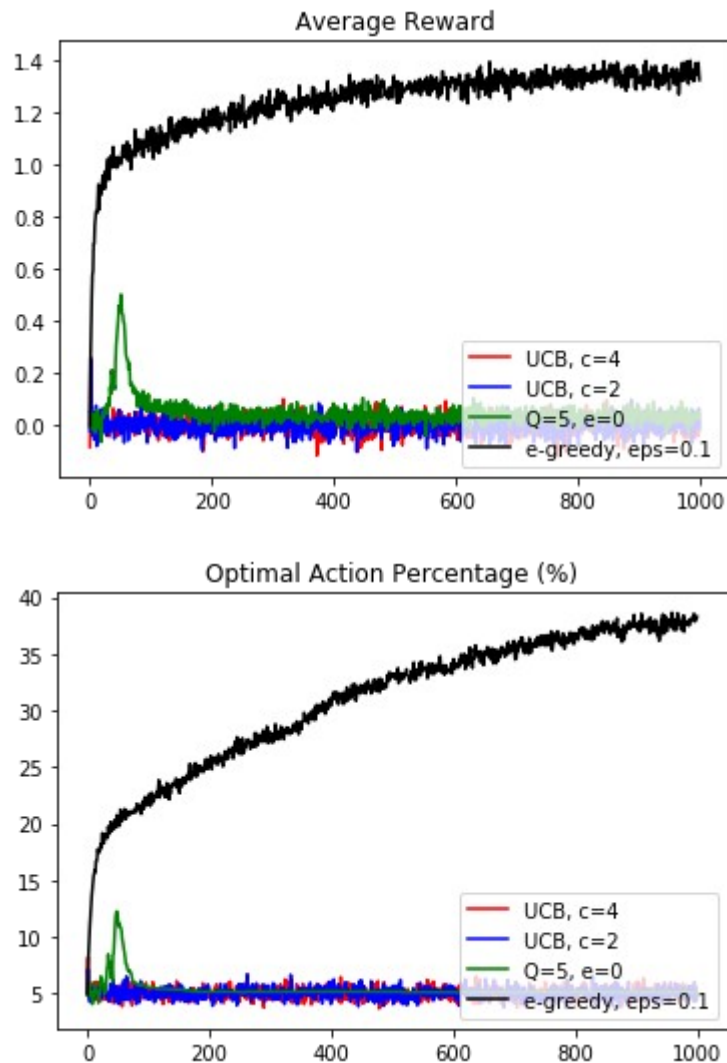




UCB for $c=4$, $c=2$, e-greedy = 0.1, $Q=5$, $e=0$; constant_step-size ; stationary



UCB for $c=4$, $c=2$, ϵ -greedy = 0.1, $Q=5$, $e=0$; constant_step-size; non-stationary



UCB performs better. For lower value of c its performance is better. Its performance is poor for constant step size.

ϵ -greedy performs best for constant step size.

Optimistic greedy has a initial increased spike but after that its performance decreases.