

Q2

Mysterious spikes

Since  $Q$  value is set to a large value of 5 compared to  $q^*$  values taken from a normal distribution with 0 mean and variance 1. When an action is selected, the reward will be less, <sup>as it is sampled from mean  $q^*$  and variance 0.01</sup> due to which the updated  $Q$  will reduce, ~~rather become negative~~ since in  $Q(s) \leftarrow Q(s) + \alpha [r - Q(s)]$ , the quantity  $r - Q(s)$  will be negative. Since  $Q$  reduces, in the next greedy pick, a different action is selected say  $a'$  and again  $Q(a')$  is reduced to lesser than  $Q(a')$ . Like this, in the first 10 picks/steps, the 10 different actions are explored and all  $Q'$  values are lesser than earlier  $Q$ s. Now, among all  $Q'$ , the max value will belong to the optimal action since it would have received highest reward when first picked, hence  $(r - Q)$  will be least negative for the optimal action. As a result, in the 11th step, the optimal action will be picked in all runs. Due to this, we see a spike in the graph of steps vs % optimal action around the 11th step.

Observation: Stationary vs Non stationary

For the stationary case, while the optimistic, greedy approach performs worse than realistic initially, once it has explored, it outperforms the realistic  $Q$ -greedy eventually. However this cannot work for non stationary cases which may require subsequent need for reexploration due to changing  $q^*$  values.  $\therefore$  In case of non stationary, the optimistic greedy approach is not as efficient.

Q1

Both graphs of average reward and % optimal action highlight the better performance of constant  $\alpha$  as compared to sample average for case of non stationary. This is because the constant  $\alpha$  gives more weight to recent rewards which is suitable for nonstationary.

Q3Unbiased constant step size Trick.1/1

$$\text{step size} = \beta_n = \frac{\alpha}{\bar{O}_n}$$

$$\bar{O}_n = \bar{O}_{n-1} + \alpha(1 - \bar{O}_{n-1}) \text{ for } n \geq 0 \text{ with } \bar{O}_0 = 0.$$

To show:  $Q_n$  is an exponential recency-weighted average without initial bias.

Proof: Basically we need to show that  $Q_n$  is independent of  $Q_1$  since all bias is included in  $Q_1$  values.

$Q$  update equation:

$$Q' = Q + \beta_n (R - Q).$$

$$Q_{n+1} = Q_n + \frac{\alpha}{\bar{O}_n} (R_n - Q_n)$$

$$Q_2 = Q_1 + \frac{\alpha}{\bar{O}_1} (R_1 - Q_1)$$

$$\bar{O}_1 = \bar{O}_0 + \alpha(1 - \bar{O}_0) = 0 + \alpha = \alpha$$

$$\frac{\alpha}{\bar{O}_1} = \frac{\alpha}{\alpha} = 1$$

$$\therefore Q_2 = Q_1 + R - Q_1 = R_1$$

$Q_2$  is independent of  $Q_1$

$$\begin{aligned} Q_3 &= Q_2 + \frac{\alpha}{\bar{O}_2} (R_2 - Q_2) \\ &= R_1 + \frac{\alpha}{\bar{O}_2} (R_2 - R_1) \end{aligned}$$

$$\bar{O}_2 = \bar{O}_1 + \alpha(1 - \bar{O}_1) = \alpha - \alpha^2 = \alpha(1 - \alpha)$$

$$Q_3 = R_1 + \frac{R_2 - R_1}{1 - \alpha}$$

$Q_3$  is independent of  $Q_1$

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

$$= Q_{n-1} [1 - \beta_n] + \beta_n R_n \quad \text{--- indep of } Q_1$$

$$Q_{n-2} + \beta_{n-2} [R_{n-2} - Q_{n-2}] = Q_{n-2} \text{ --- indep of } Q_1$$

Since  $Q_n$  depends on  $Q_{n-1} \dots$  which only depends on  $Q_2$  which is independent of  $Q_1$ ,  $Q_n$  is unbiased



Q4

Comparing UCB, greedy optimistic and realistic  $\epsilon$ -greedy,  
for stationary, greedy optimistic starts out as the 1/1

worst performing method but later outperforms both UCB and realistic methods. Initial exploration causes the low initial average rewards, however later, it performs the best as it greedily chooses the best action each time which is suitable in a stationary setting. Among UCB and realistic  $\epsilon$ -greedy, UCB performs marginally better. This could be due to adaptive exploration in UCB which increases  $\beta$  value more for the actions not selected in turn encouraging exploration.

For the non stationary case, UCB performs poorly and in the long run even the optimistic greedy approach fails to perform well due to not exploring at the later times.  $\epsilon$ -greedy performs similar to UCB and is perhaps easier to extend to non stationary cases as compared to UCB.