

## 1. 1 point. (RL2e 2.2) Exploration vs. exploitation.

**Written:** Consider a  $k$ -armed bandit problem with  $k = 4$  actions, denoted 1, 2, 3, and 4. Consider applying this problem to a bandit algorithm using  $\epsilon$ -greedy action selection, sample-average action-value estimates, and initial estimates of  $Q_1(a) = 0$ , for all  $a$ . Suppose the initial sequence of actions and rewards is  $A_1 = 1, R_1 = -1, A_2 = 2, R_2 = 1, A_3 = 2, R_3 = -2, A_4 = 2, R_4 = 2, A_5 = 3, R_5 = 0$ . On some of these time steps did this definitely occur? On which time steps did this possibly have occurred?

Timestep	$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$
0	0	0	0	0
1	-1	0	0	0
2	-1	1	0	0
3	-1	-1.5	0	0
4	-1	-0.5	0	2
5	-1	0.5	0	2

Action selected at Random  
 $A_1 T=1, A_1=1 \rightarrow R=-1$ , therefore Action must have been random  
At  $T=2$ ,  $A_2=2, R=1$   
Most second greedy  
At  $T=3$ ,  $A_3=2, ER=-0.5$   
Selecting must be random as Reward = 0 is good  
At  $T=4$ ,  $A_4=4, R=2$   
Selecting Greedy  
At  $T=5$ ,  $A_5=3, R=0$   
Selection must have been random as maximum reward not chosen

## 2. 1 point. (RL2e 2.4) Varying step-size weights.

**Written:** If the step-size parameters,  $\alpha_n$ , are not constant, then the estimate  $Q_n$  is a weighted average of previously received rewards with a weighting different from that given by Equation 2.6. What is the weighting on each prior reward for the general case, analogous to Equation 2.6, in terms of the sequence of step-size parameters?

$$\begin{aligned} \theta_{n+1} &= \theta_n + \alpha_n [R_n - \theta_n] && \text{-(2.5)} \\ &= \alpha_n R_n + (1-\alpha_n) \theta_n && \text{expanding } \theta_n \text{ from eqn (2.5)} \\ &= \alpha_n R_n + (1-\alpha_n) [\theta_{n-1} + \alpha_{n-1} [R_{n-1} - \theta_{n-1}]] \\ &= \alpha_n R_n + \alpha_{n-1} R_{n-1} - \alpha_n \alpha_{n-1} R_{n-1} + (1-\alpha_n)(1-\alpha_{n-1}) \theta_{n-1} \\ &= \alpha_n R_n + (1-\alpha_n) \alpha_{n-1} R_{n-1} + (1-\alpha_n)(1-\alpha_{n-1}) \theta_{n-2} \\ &\vdots && \\ \theta_{n+1} &= \sum_i^n \alpha_i R_i - R_{n-1} \prod_i^n \alpha_i + \theta_1 \prod_i^n (1-\alpha_i) \end{aligned}$$

Q2

## 2 points. Bias in Q-value estimates.

This question is required for 5180 and extra credit for 4180

**Written:** Recall that  $Q_n \triangleq R_{n-1} + \dots + R_1$  is an estimate of the true expected reward  $q_*$  of an arbitrary arm  $a_{n-1}$ . We say that an estimate is biased if the expected value of the estimate does not match the true value, i.e.,  $E[Q_n] \neq q_*$  (otherwise, it is unbiased).

(a) Consider the sample-average estimate in Equation 2.1. Is it biased or unbiased? Justify your answer with brief words and equations.

$$Q_n = \frac{1}{n-1} (R_1 + R_2 + \dots + R_{n-1}) \quad \text{-(eq 2.1)}$$

$$\text{as we know } E[R_n] = q_*$$

$$\begin{aligned} E[Q_n] &= \frac{1}{n-1} E[R_1 + R_2 + \dots + R_{n-1}] \\ &= \frac{1}{n-1} [E[R_1] + E[R_2] + \dots + E[R_{n-1}]] \\ &= \frac{1}{n-1} \times (n-1) q_* \\ &= q_* \end{aligned}$$

Therefore Sample average Estimate in eq 2.1 is unbiased

For the remainder of the question, consider the exponential recency-weighted average estimate in Equation 2.5.

Assume that  $0 < \alpha < 1$  (i.e., it is strictly less than 1).

(b) If  $Q_0 = 0$ , is  $Q_n$  (for  $n > 1$ ) biased? Justify your answer with brief words and equations.

$$\begin{aligned} \theta_{n+1} &= (1-\alpha)^n \theta_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i \\ \theta_{n+1} &= \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i \quad \text{as } \theta_1 = 0 \end{aligned}$$

Expectation  $E[\theta_{n+1}]$

$$E[\theta_{n+1}] = E\left[\sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i\right]$$

$$= \sum_{i=1}^n \alpha (1-\alpha)^{n-i} E[R_i]$$

$$= \sum_{i=1}^n \alpha (1-\alpha)^{n-i} q_*$$

$$\left\{ \begin{array}{l} \text{When } \sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1 \\ E[\theta_{n+1}] = q_* , \text{ it will be unbiased} \end{array} \right.$$

else biased

(c) Derive condition(s) for  $Q_1$  for when  $Q_n$  will be unbiased.

as per part b

$$\begin{aligned} E[\theta_{n+1}] &= E\left[(1-\alpha)^n \theta_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i\right] \\ &= (1-\alpha)^n E[\theta_1] + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} E[R_i] \\ &= (1-\alpha)^n E[\theta_1] + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} q_* \end{aligned}$$

It will be Unbiased when

$$\theta_1 = 0 \quad \& \quad \sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1$$

Q3

In the long run  $\epsilon = 0.01$  perform better

$\epsilon = 0$ , will always choose the action which discarded first & being greedy

$\epsilon = 0.1$  = optimal action  $\theta_0$

$$0.9 \times 1 + 0.1 \times \frac{1}{10} = 91\%$$

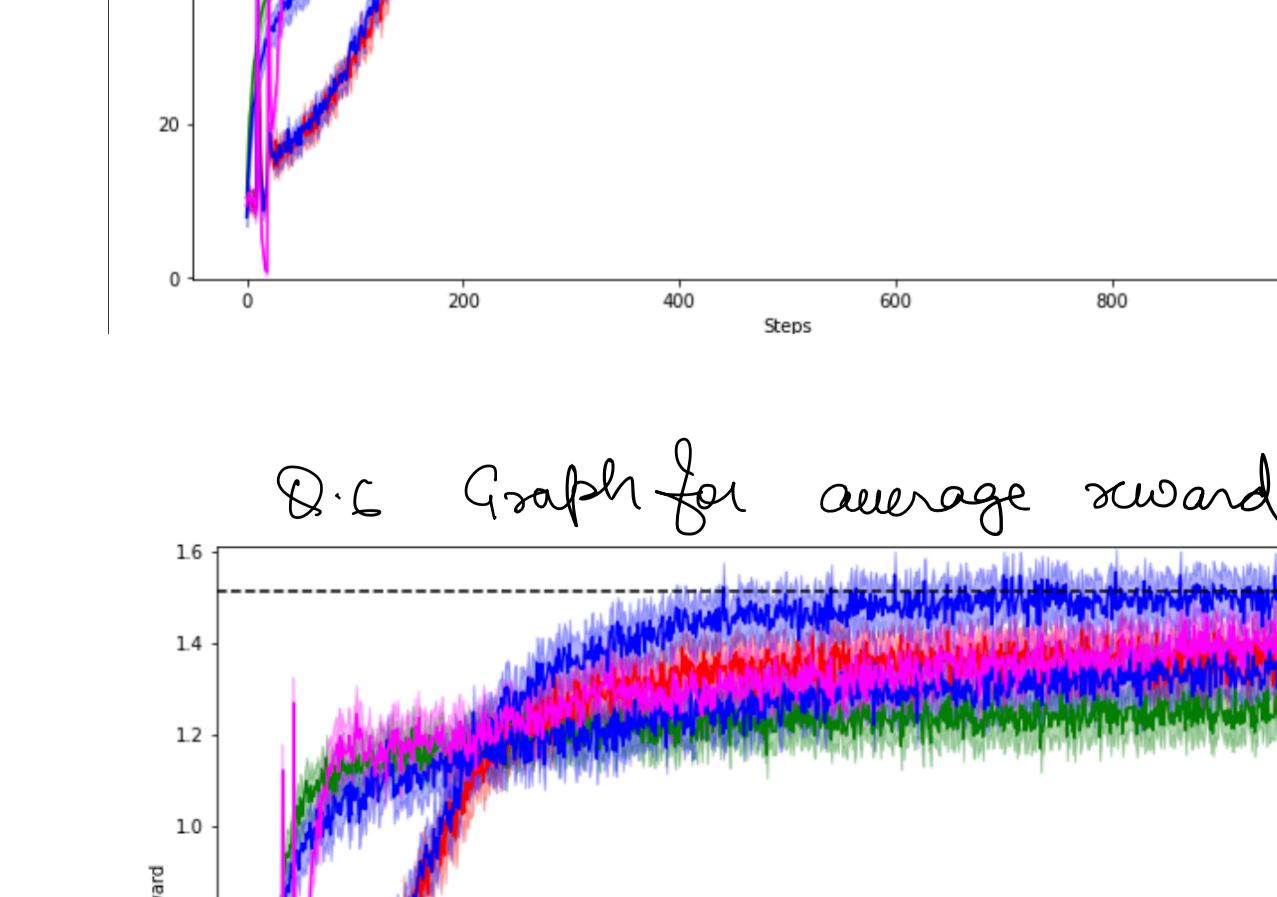
$\epsilon = 0.01$  = optimal action  $\theta_0$

$$0.99 + 0.01 \times \frac{1}{10} = 99.1\%$$

Although greedy has 100% chances of selecting best action due to lack of exploration it will get stuck

$\epsilon = 0.1 \leftarrow$  has more reward than  $\epsilon = 0.01$  initially but after some time when explored enough, it will also become more exploitative than exploratory

$\therefore \epsilon = 0.01$  cumulatively will return higher reward

Q4

we see the spike for the optimistic initialisation in class.

Spike appears but not large enough as  $n \rightarrow \infty$ , asymptotically it's possible

Q4. Graph Optimal action



Plot for Q4 optimal action

Q5

we see the spike for the optimistic initialisation in class.

Spike appears but not large enough as  $n \rightarrow \infty$ , asymptotically it's possible

Q5. Graph average reward



Cases with  $\epsilon$  reach asymptotic & with initial Value has the spike