# Q1)

$A_1 = 1 \qquad R_1 = -1$

$A_2 = 2 \qquad R_2 = 1$

$A_3 = 2 \qquad R_3 = -2$

$A_4 = 2 \qquad R_4 = 2$

$A_5 = 3 \qquad R_5 = 0$

$Q_0 =$

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

$Q_1 =$

| -1 | 0 | 0 | 0 | 0 |
|----|---|---|---|---|

$Q_2 =$

| -1 | 1 | 0 | 0 | 0 |
|----|---|---|---|---|

$Q_3 =$

| -1 | -1 | 0 | 0 | 0 |
|----|----|---|---|---|

$Q_4 =$

| -1 | 1 | 0 | 0 | 0 |
|----|---|---|---|---|

$Q_5 =$

| -1 | 1 | 0 | 0 | 0 |
|----|---|---|---|---|

→ There is a possibility that action $A_3$ could have been a random selection

→ We can ~~def~~ positively state that, actions ~~A~~ $A_4$ & $A_5$ are definitely $\varepsilon$-cases where selection was done ~~so~~ randomly to explore.

(2)

$$Q_{m+1} = Q_m + \alpha_m \left[ R_m - Q_m \right]$$

$$Q_{m+1} = \alpha_m R_m + (1-\alpha_m) Q_m$$

$$\cancel{Q_{m+1} = \alpha_m R_m + (1-\alpha)\left[ a_{m-1} \cancel{+} \alpha \right.}$$

$$Q_{m+1} = \alpha_m R_m + (1-\alpha_m)\left[ \alpha_{m-1} R_{m-1} + (1-\alpha_{m-1}) Q_{m-1} \right]$$

$$Q_{m+1} = \alpha_m R_m + \alpha_{m-1}(1-\alpha_m) R_{m-1} + (1-\alpha_m)(1-\alpha_{m-1})\left[ \alpha_{m-2} R_{m-2} + (1-\alpha_{m-2}) Q_{m-2} \right]$$

$$Q_{m+1} = \alpha_m R_m + (1-\alpha_m)\alpha_{m-1} R_{m-1} + \ldots + (1-\alpha_m)(1-\alpha_{m-1})\ldots(1-\alpha_2)\alpha_1 R_1$$
$$+ (1-\alpha_m)(1-\alpha_{m-1})\ldots(1-\alpha_1) Q_1$$

$$Q_{m+1} = \left[ \prod_{i=1}^{n} (1-\alpha_i) \right] Q_1 + \sum_{i=1}^{n}\left[ \alpha_i R_i \cdot \prod_{j=i+1}^{m} (1-\alpha_j) \right]$$

# 3.

## (a)

Equation 2.1:

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i = a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i = a}}$$

By looking at the equation, we can say that all the rewards are equally weighted and therefore this equation gives the Expectation $(E)$ of $Q_t(a)$.

By the law of large numbers, $Q_t(a)$ will eventually converge to $q_*^*(a)$ as the number of samples reach infinity. So we can say that the estimate is unbiased if a sufficiently large number is chosen.

$$E(Q_n) = E\left[\frac{R_1 + R_2 + \cdots R_{m-1}}{n-1}\right] = E[E(R_{t-1})] = E[q^*] = q^*$$

$$E(Q_n) = q^*$$

23)

(d)

$$Q_n = Q_{n-1} + \alpha \left[ R_{n-1} - Q_{n-1} \right]$$

If $Q_1 = 0$ and is free from $q^*$, initially for $n > 1$,

$Q_n$ will the chased

$$Q_2 = Q_1^0 + \alpha \left[ R_1 - Q_1^0 \right] = \alpha R_1$$

(c)

For $Q_n$ to be unbiased, $Q_1$ should be set to $q^*$

$$Q_1 = q^*$$

3.

(d)

$$Q_{m+1} = Q_m + \alpha \left[ R_m - Q_m \right]$$

$$E(Q_{m+1}) = E(Q_m) + \alpha \left[ E(R_m) - E(Q_m) \right]$$

m. d. l,

at $m \to \infty$, $E(R_m) \to q^*$

$$E(Q_{m+1}) = E(Q_m) + \alpha \left[ q^* - E(Q_m) \right]$$

$$E(Q_{m+1}) = (1-\alpha) E(Q_m) + \alpha q^*$$

$$= (1-\alpha) \left\{ (1-\alpha) E(Q_{m-1}) + \alpha q^* \right\} + \alpha q^*$$

$$= (1-\alpha)^m E(Q_1) + \alpha q^* + \alpha (1-\alpha) q^* + \alpha (1-\alpha)^2 q^* + \cdots$$

$$= (1-\alpha)^m E(Q_1) + \sum_{i=1}^{m-1} \alpha (1-\alpha)^i q^*$$

$$= (1-\alpha)^m E(Q_1) + \alpha q^* \sum_{i=1}^{m-1} (1-\alpha)^i$$

$$\longrightarrow \alpha q^* \left\{ \frac{1 - (1-\alpha)^{m-1}}{1 - (1-\alpha)} \right\}$$

$$= q^* \left\{ 1 - (1-\alpha)^{m-1} \right\}$$

$$E(Q_{m+1}) = (1-\alpha)^m E(Q_1) + q^* \{1 - (1-\alpha)^{m-1}\}$$

as $m \to \infty;$

$$E(Q_{m+1}) = \underbrace{(1-\alpha)^m}_{0} Q_1 + q^* [1 - \underbrace{(1-\alpha)^{m-1}}_{0}]$$

$$\boxed{E(Q_{m+1}) = q^*}$$

30)

Exponential recency - weighted average is ~~of the form~~ given by

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

It is of the form,

New Estimate = Old Estimate + Step - Size [ Target - Old Estimate ]

The nature of the equation is to decrease the difference

between the Reward and its estimate

4.

$$\Pi_t(a) = \frac{e^{H_t(a)}}{e^{H_t(a)} + e^{H_t(b)}}$$

$$\boxed{\Pi_t(a) = \frac{1}{1 + e^{H_t(b) - H_t(a)}}}$$

Sigmoid :

$$\boxed{f(x) = \frac{1}{1 + e^{-x}}}$$

5) The algorithm ~~makes~~ ensures that the probability of its optimal action being selected is a but ones $1 - \varepsilon$.

It is a but ones $1 - \varepsilon$ because there is a chance that the optimal action will be chosen when the ~~also~~ algorithm tries to explore ~~make~~ by making random selection.

6)

UCB starts slow since it starts off exploring all its actions. The sharp increase comes from the algorithm repeatedly choosing the action which it has the highest Q-value ensuring that it takes the most optimal action.

$$A_t = \arg\max_a \left[ Q_t(a) + c \sqrt{\frac{\ln(t)}{N_t(a)}} \right]$$

As the UCB keeps choosing the optimal actions, $\sqrt{\frac{\ln(t)}{N_t(a)}}$ decreases as $N_t(a)$ increases faster that $\ln(t)$. This forces UCB to choose other actions until the percacme of the optimal action becomes the heighest. This is the reason for the sudden decrease in the % of optimal action chosen.