

Luke Davidson

CS 5180

Ex. 1

Q_t(a)					A_t	R_t
t	1	2	3	4		
1	0	0	0	0	1	-1
2	-1	0	0	0	2	1
3	-1	1	0	0	2	-2
4	-1	-0.5	0	0	2	2
5	-1	0.33	0	0	3	0

$$Q_4(2) = \frac{1-2}{2} = -\frac{1}{2}$$

$$Q_5(2) = \frac{1-2+2}{3} = \frac{1}{3} = 0.33$$

t=1: May have occurred. Since all Q_t(a) are 0, both the ε and 1-ε options would select arbitrarily, so we don't know which it was.

t=2: May have occurred. Similar to t=1, Q_2({2,3,4}) are all the max and = 0. For ε case, a random selection of 1, 2, 3 or 4 would be made, so 2 would be possible. For 1-ε case, a random selection of 2, 3 or 4 would be made, so 2 also possible.

t=3: May have occurred. For ε case, a random selection of 1-4 would have been made, so 2 possible. For 1-ε case, 2 would be selected since it is the max Q_3(a).

t=4: Definitely occurred. 1-ε case would have selected 3 or 4 since the Q_4({3,4}) are maxes. Since 2 was selected, we know it was by random, or ε case.

t=5: Definitely occurred. 1-ε case would have selected 2 since 0.33 is max Q_5(a). 3 was selected so we know it was by random, or ε case.

$$\begin{aligned}
 2) Q_{n+1} &= Q_n + \alpha_n [R_n - Q_n] \\
 &= Q_n + \alpha_n R_n - \alpha_n Q_n \\
 &= \alpha_n R_n + \underbrace{\alpha_n(1-\alpha_n)}_{\text{factored}} \\
 &= \alpha_n R_n + (1-\alpha_n) \left(Q_{n-1} + \alpha_{n-1} [R_{n-1} - Q_{n-1}] \right) \\
 &= \alpha_n R_n + (1-\alpha_n) \left(R_{n-1} \alpha_{n-1} + (1-\alpha_{n-1})(Q_{n-1}) \right) \\
 &= \underbrace{\alpha_n R_n + (1-\alpha_n)(R_{n-1} \alpha_{n-1})}_{\text{factored}} + \underbrace{(1-\alpha_n)(1-\alpha_{n-1})(Q_{n-1})}_{\text{factored}} \\
 &= (1) \cdot (Q_{n-2} + \alpha_{n-2} [R_{n-2} - Q_{n-2}])
 \end{aligned}$$

$$\begin{aligned}
 &= \alpha_n R_n + (1-\alpha_n)(R_{n-1} \alpha_{n-1}) + (1-\alpha_n)(1-\alpha_{n-1})(R_{n-2} \alpha_{n-2}) \\
 &\quad + (1-\alpha_n)(1-\alpha_{n-1})(1-\alpha_{n-2})(Q_{n-2})
 \end{aligned}$$

eventually will lead to

$$Q_{n+1} = \left[\prod_{i=1}^n (1-\alpha_i) \right] Q_1 + \sum_{i=1}^{n-1} \alpha_i \prod_{i=1}^n (1-\alpha_{i+1}) b_i + \alpha_n R_n$$

$$3) \text{a) } \text{eq. 2.1} = \frac{\sum_{i=1}^{t-1} R_i \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}} = E[Q_n]$$

$$z_a = E[R_i | A_i=a]$$

$$E[Q_n] = \dots \Rightarrow E\left[\frac{R_1 + \dots + R_{n-1}}{n-1}\right]$$

$\hookrightarrow Q_n = \frac{R_1 + \dots + R_{n-1}}{n-1}$

$$= E\left[\frac{\sum_{i=1}^{n-1} R_i}{\sum_{i=1}^{n-1} \mathbb{1}_{A_i=a}}\right]$$

for eq. 2.1

$$E\left[\frac{\sum_{i=1}^{t-1} R_i \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}\right]$$

Comparing both, they are equivalent through equating n and t. so it is unbiased.

$$3b) \text{ Q}_n \stackrel{?}{=} Q_n + \alpha [R_n - Q_n]$$

$$Q_{n+1} \Rightarrow (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

$$\text{for } Q_1 = 0$$

$$Q_{n+1} = \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

$$Q_{n+1} = \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

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

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

} this statement is not true, so when $Q_1 = 0$,
the exponential recency-weighted average
estimate is baised.

3c.) for Q_n to be unbiased

$$E[Q_n] = q_* \quad \text{for } Q_n = (1-\alpha)^{n-1} Q_1 + \sum_{i=1}^{n-1} \alpha (1-\alpha)^{n-i} R_i$$

$$\rightarrow E[Q_n] = Q_1 E[(1-\alpha)^{n-1}] + \sum_{i=1}^{n-1} E[R_i] \alpha (1-\alpha)^{n-i}$$

must solve for α Q_1 , or rule for Q_1 , where

$$E[Q_n] = E\left\{\frac{\sum R_i}{n-1}\right\} = q_* = E[R_t | A_t = a]$$

$$Q_1 = \frac{E\left[\frac{\sum R_i}{n-1}\right] - \sum_{i=1}^{n-1} E[R_i] \alpha (1-\alpha)^{n-i}}{E[(1-\alpha)^{n-1}]}$$

$n \rightarrow \infty$

$$3d) Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

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

as $n \rightarrow \infty$,

$$(1-\alpha)^\infty Q_1 + \sum_{i=1}^{\infty} \alpha (1-\alpha)^{\infty-i} R_i \rightarrow \textcircled{D} \quad b/c \quad 0 < \alpha < 1,$$

so $(1-\alpha)^\infty$ will continue to get smaller and smaller and eventually get to 0.

$$E[Q_n] = E\left[\frac{\sum_{i=1}^n R_i}{n-1}\right] \text{ as } n \rightarrow \infty$$

also $\rightarrow \textcircled{D}$ b/c the denom will be ∞ . Therefore

as $n \rightarrow \infty$ Q_n is unbiased since both = 0.

3e). We should expect the exponential recency-weighted average will be biased in practice due to the original effect Q_1 has on it. For

$$Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i,$$

as n increases, $(1-\alpha)^n$ will become smaller and smaller (given $0 < \alpha < 1$), although never 0. Because it will never achieve the value of 0, there will always be a diminishing bias towards Q_1 , thus it will behave as a bias estimator.

$$4) \Pr\{A_t = a\} = \frac{e^{H_t(a)}}{\sum_{b=1}^L e^{H_t(b)}} = \pi_t(a)$$

logistic

$$\frac{L}{1 + e^{-k(x-x_0)}}$$

$$2 \text{ actions} \quad \Pr\{A_t = a\} = \frac{e^{H_t(a)}}{e^{H_t(a)} + e^{H_t(b)}} \cdot \frac{\frac{1}{e^{H_t(a)}}}{\frac{1}{e^{H_t(b)}}}$$

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

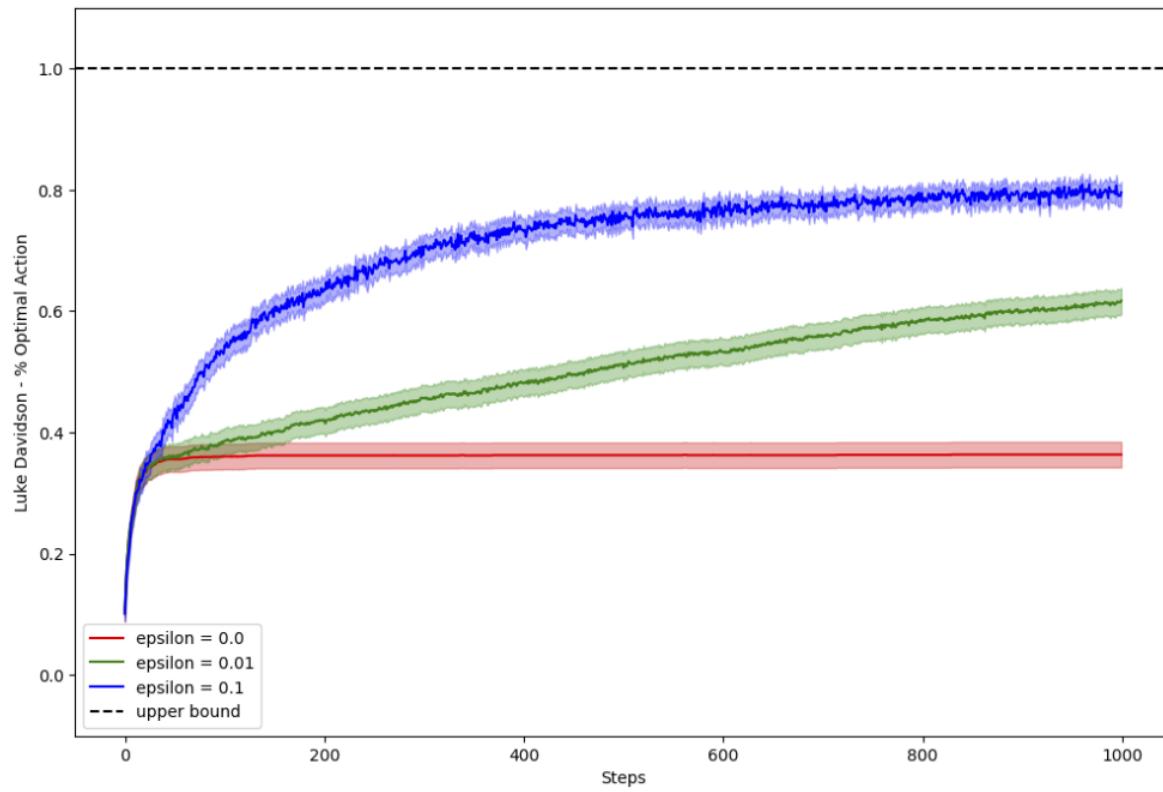
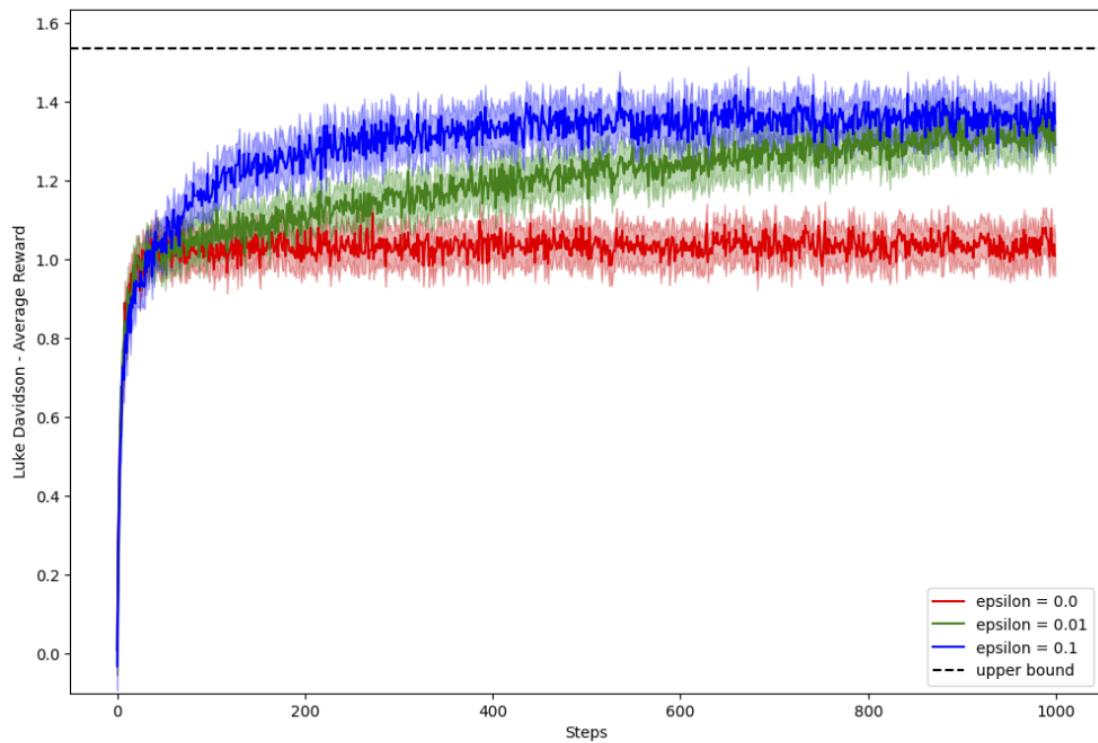
which is in same form as

logistic

$$\boxed{\frac{L}{1 + e^{-k(x-x_0)}}}$$

where $L=1$
and
 $-k(x-x_0) = H_t(b) - H_t(a)$

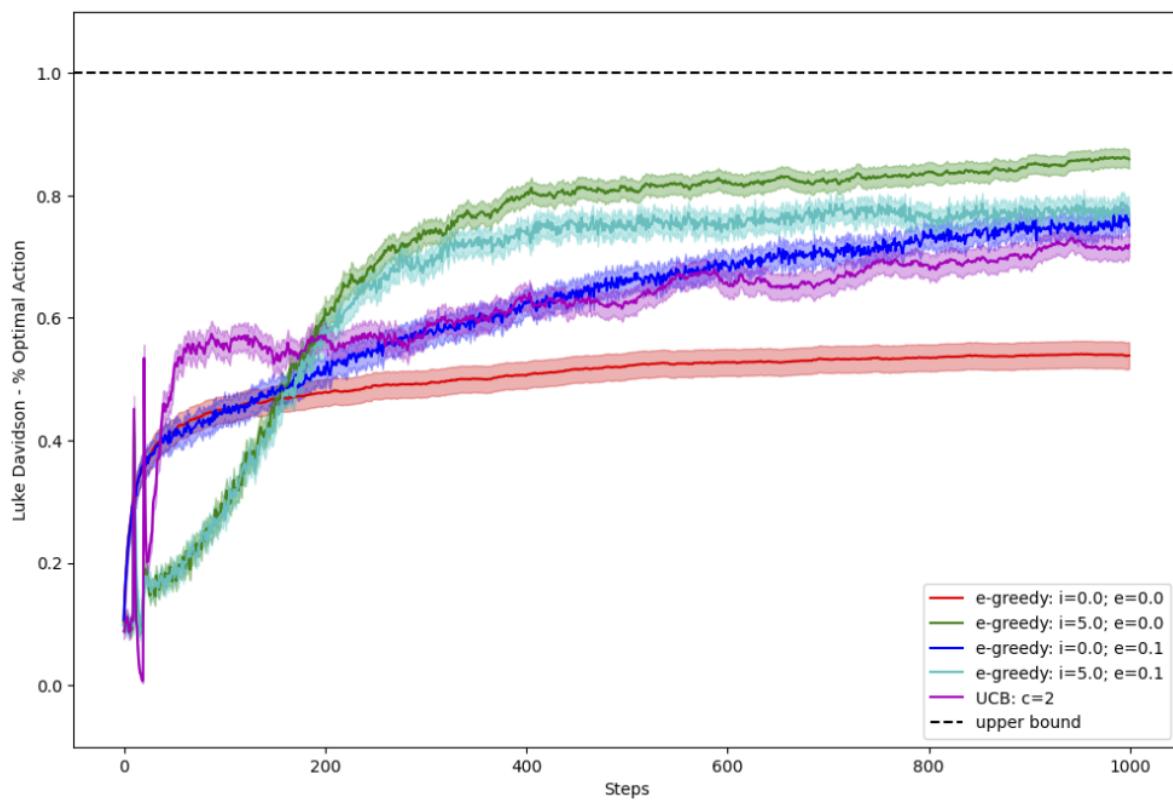
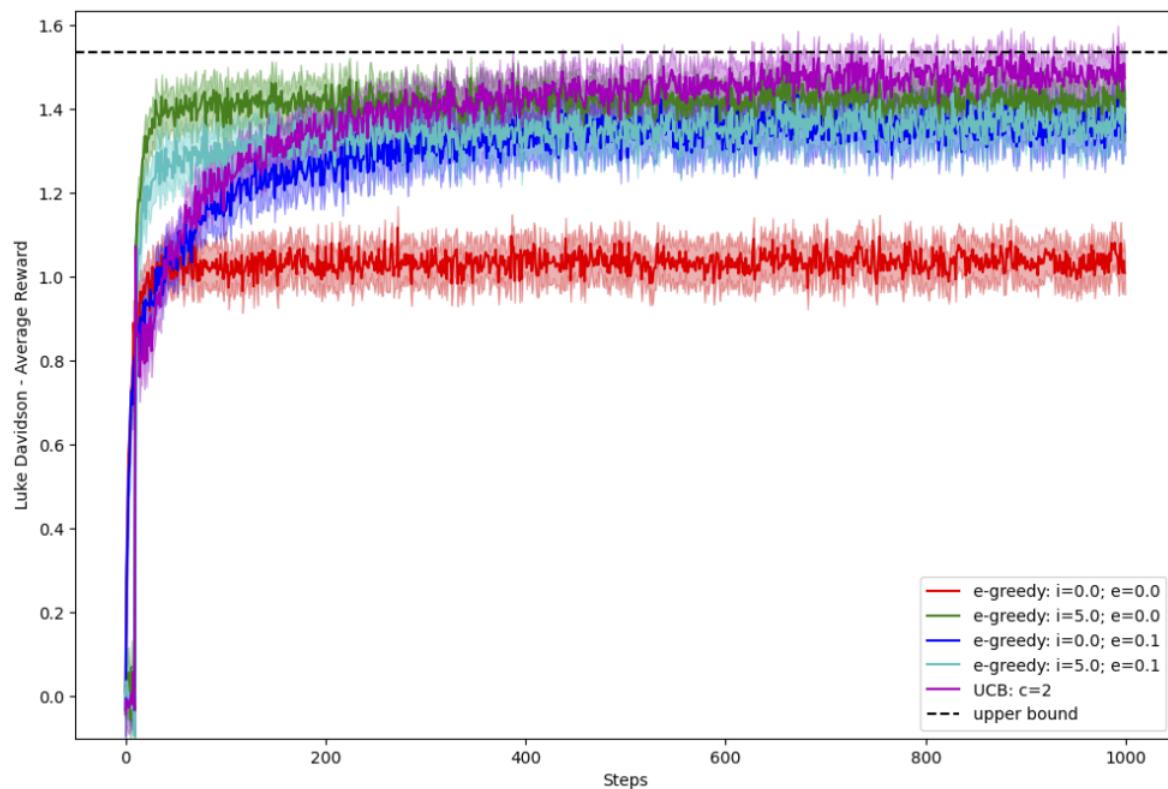
5 Plot



5.) Written: The average rewards that the plots converge to essentially represent the upper limit each of those individual policies would be able to reach.^(in theory) As the # of steps increase, eventually to infinity, each of these policies will have different levels of optimality, represented by the values they converge to. At some point, each policy will not be able to gain more rewards per action.

This is due to the nature of the policies, which itself is due to the hyperparameters. These policies behave differently with different epsilon values. A policy with a high epsilon value will more often times land on a random action choice compared to lower epsilon values, thus "exploring" more and "exploiting" less. A perfect ratio of exploring and exploiting is desired to obtain a higher optimal convergence. Too much exploration and the agent will not exploit its knowledge of which action may be best enough. Too much exploitation and the agent will not learn the environment well enough. This is seen in our plots with $\epsilon = \{0, 0.01, 0.1\}$.

6 Plot



(v) Written: At the beginning rewards are low because all actions are treated as maxes due to $N_t(a) = 0$, so they are selected randomly. With a higher " c ", the variance in the estimate of a 's value is amplified. After all the actions are selected an initial time, eventually a high reward action will be selected randomly, which will greatly increase the average reward of that step, resulting in the spike seen. After the first spike, and similarly to the first part, c aids to not allow the agent to act in a greedy way, so it will likely not select the highest reward action again. This results in the decrease after the initial spike. As t grows, This process more or less so continues until ...