

Miles Vollmer
RL Homework 2 EX-1

Q1:

Work:

	T	#1	R1	#2	R2	#3	R3	#4	R4
greedy = rand	1	1	-1						
g = rand(2,3,4)	2			1	1				
g = 2	3			2	-2				
g = rand(3,4)	4			3	2				
g = 2	5					1			

Explanation:

There was definitely an epsilon selection at timesteps 4 and 5. I know this because at timestep 1 both greedy and epsilon would just choose randomly, so can't guarantee an epsilon selection at this point. At timestep 2 greedy would choose randomly between actions 2, 3, and 4 and 2 is selected so can't guarantee an epsilon selection at this point. At timestep 3 greedy would choose 2 and 2 is selected so again can't guarantee an epsilon selection at this point. At timestep 4 greedy would randomly select between 3 and 4, as the expected reward for 1 and 2 are both negative at this point. However, 2 was selected at this point so we can guarantee that this was an epsilon selection. At timestep 5 greedy would select action 2 as its expected reward is now higher than any other action. However, 3 was selected at this point meaning that we can guarantee this was an epsilon selection.

Q2.

Work:

$$\begin{aligned}
 2 \quad Q_{n+1} &= Q_n + \alpha_n [R_n - Q_n] \\
 &= \alpha_n R_n + (1 - \alpha_n) Q_n \\
 &= \alpha_n R_n + (1 - \alpha_n) [Q_{n-1} + \alpha_{n-1} [R_{n-1} - Q_{n-1}]] \\
 &= \alpha_n R_n + Q_{n-1} + \alpha_{n-1} R_{n-1} - \alpha_{n-1} Q_{n-1} \\
 &= \alpha_n R_n + \alpha_{n-1} R_{n-1} + (1 - \alpha_{n-1}) Q_{n-1} \\
 &= \alpha_n R_n + (1 - \alpha_n) (\alpha_{n-1} R_{n-1} + (1 - \alpha_{n-1}) Q_{n-1}) \\
 &= \alpha_n R_n + \alpha_{n-1} R_{n-1} - \alpha_n \alpha_{n-1} R_{n-1} + (1 - \alpha_n) (1 - \alpha_{n-1}) Q_{n-1} \\
 &\quad \downarrow \\
 &= \sum_i \alpha_i R_i - R_{n-1} \prod_i \alpha_i + Q_1 \prod_i (1 - \alpha_i)
 \end{aligned}$$

Explanation: If we take our equation 2.5 and break it down similarly to how it was done for 2.6 however in this case we don't consider alpha to be constant. This results in a weighting similar to that of 2.6 with some key differences (i.e. The new product portions)

Q4.

Work:

4. If we have 2 actions then
$$\Pr(A_+) = \frac{e^{H_+(1)}}{e^{H_+(1)} + e^{H_+(2)}}$$

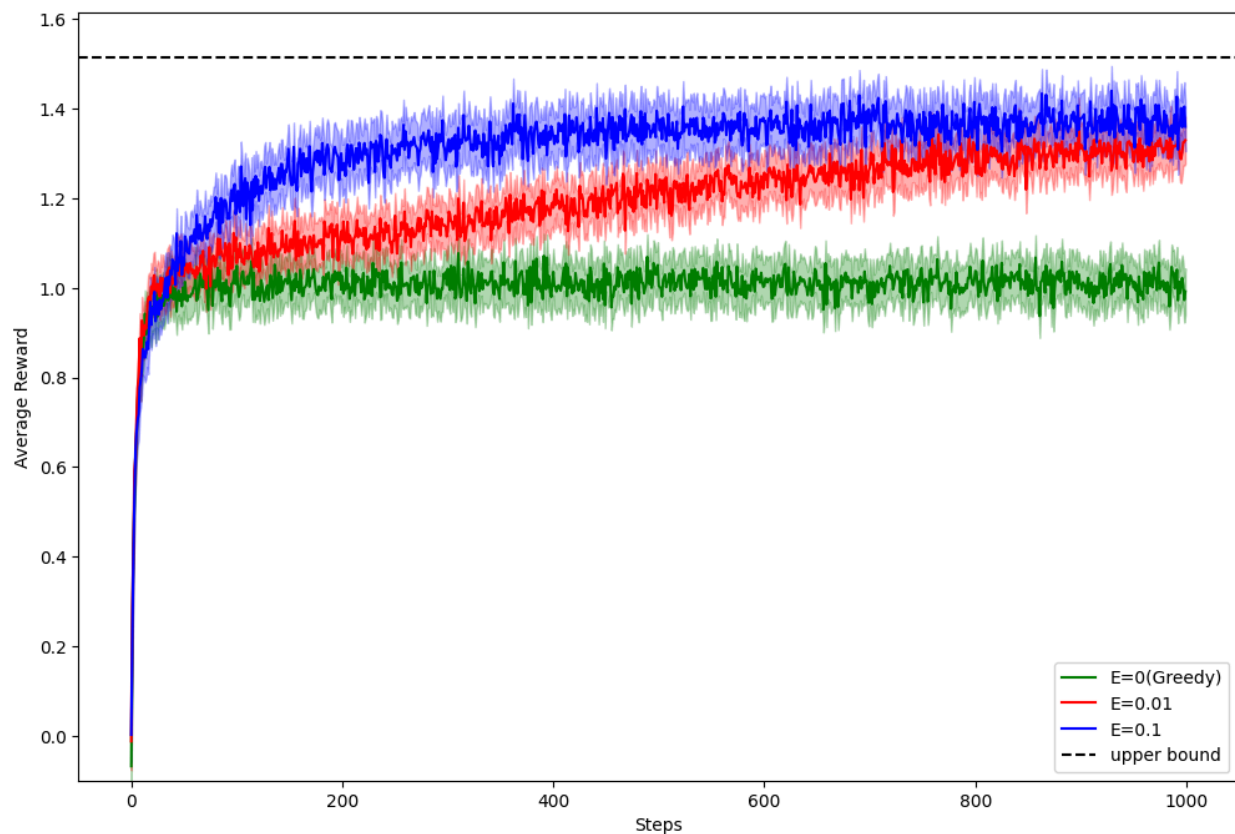
Let's write as $\frac{e^a}{e^a + e^b}$ where $a = H_+(1)$ $b = H_+(2)$
now let's multiply by $\frac{e^c}{e^c}$ where $c = -b$, we get
$$\frac{e^{a+c}}{e^{a+c} + e^{b+c}} = \frac{e^{a-b}}{e^{a-b} + e^{b-b}} = \frac{e^{a-b}}{e^{a-b} + 1}$$

you can write $\text{sigmoid}(x)$ as $\frac{e^x}{e^x + 1}$
so if $x = a - b$ then $\text{sigmoid}(x) = \frac{e^{a-b}}{e^{a-b} + 1}$
which proves in the case of two actions the softmax distribution is same as the sigmoid.

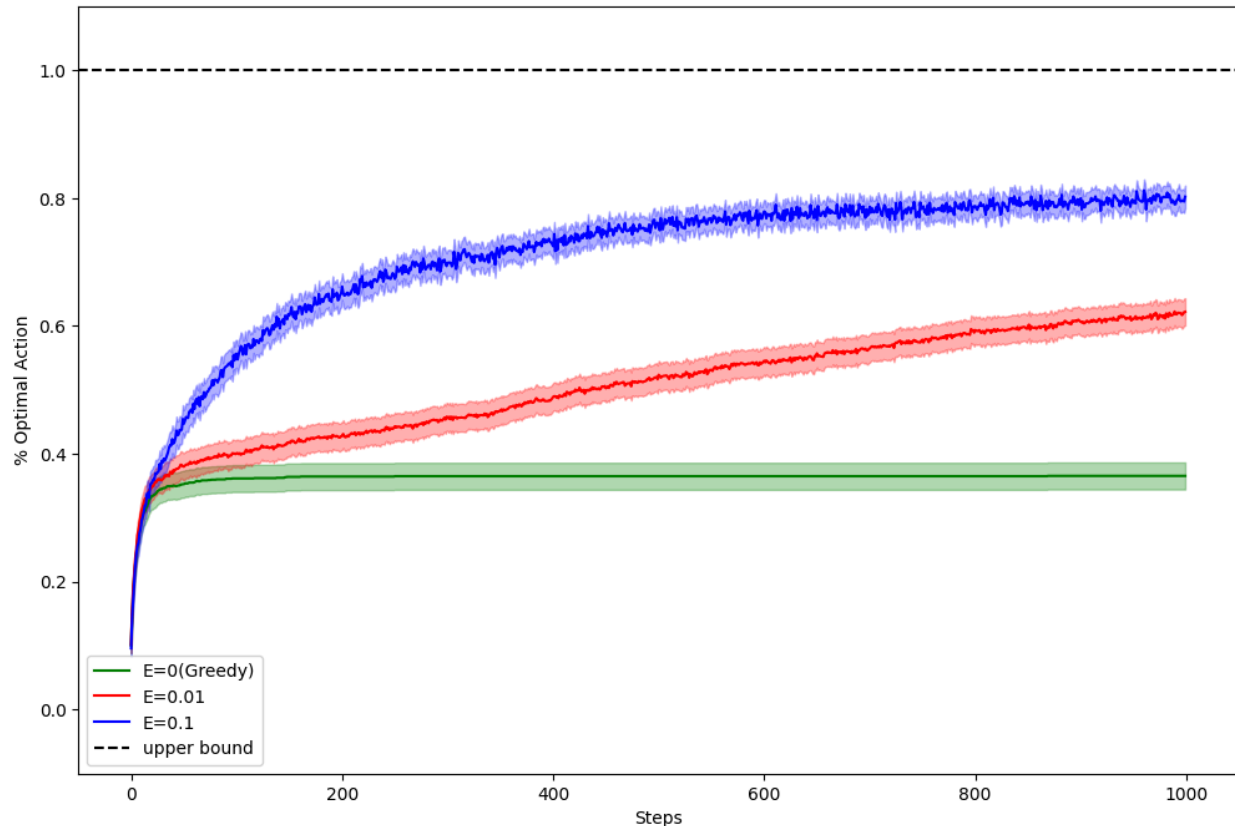
Explanation: see picture above

Q5.

Plot:



Average Reward of E-Greedy with different epsilon



Optimal Action Percentage of E-Greedy with different epsilon

Written: E=0 converges around 1, E=0.01 converges around 1.25, and E=0.1 converges around 1.35. The algorithms converge to different values due to their different epsilon values because these epsilon values dictate how much exploration the algorithm performs. So in this graph we see that the E=0.1 algorithm performs best over the 1000 steps. This is because the E=0 does no exploration and only chooses what it thinks is the best action, so it never learns the optimal and converges at 1. The E=0.01 algorithm is interesting because it performs worse over these 1000 time steps as it explores slower than E=0.1 as it is selecting what it thinks is the best option 99% of the time. However with E=0.01 once it learns the optimal action it picks it 9% more often than with E=0.1 so if we extended this graph we would see E=0.01 take over as the best performing algorithm.

Q6.
Plot:

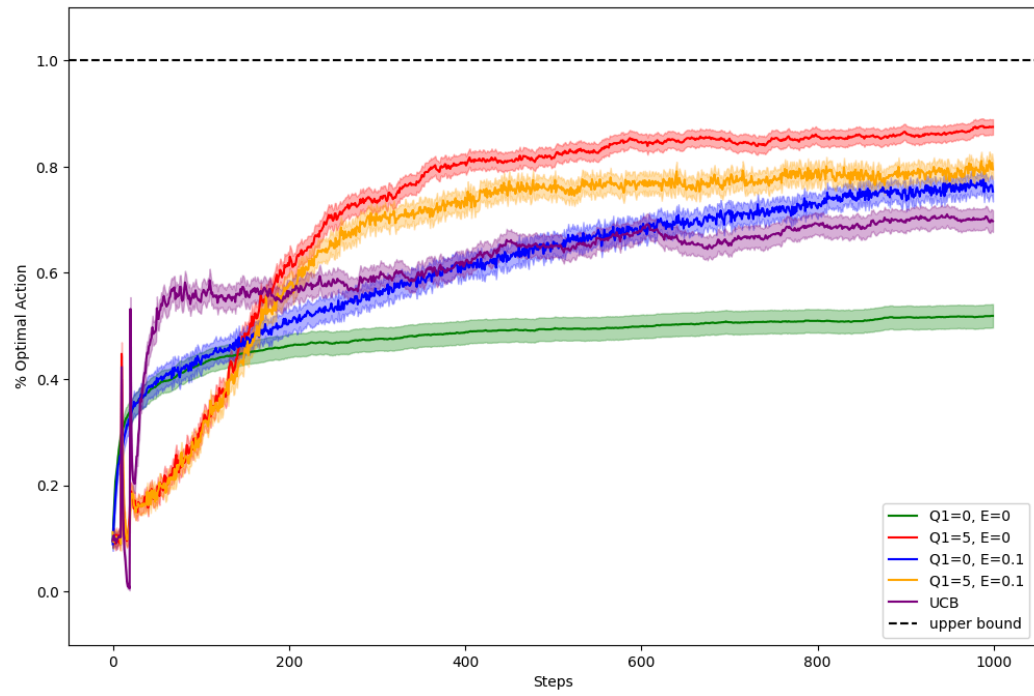


Fig 2.3: Reproduced

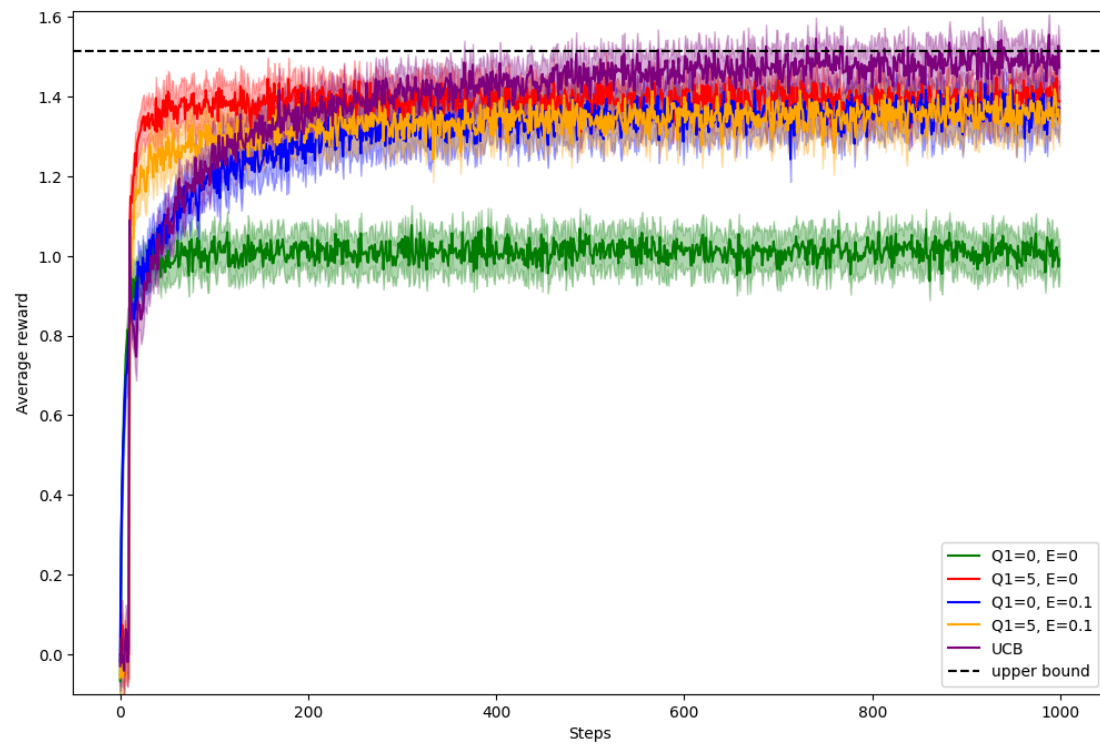


Fig 2.4: Reproduced

Written:

The spikes that appear in the UCB algorithms at the beginning of the two figures are due to how the UCB algorithm goes about its initial exploration. In the UCB algorithm, the initial surge in performance arises from the algorithm's strategy of favoring exploration, particularly for actions that have been chosen less frequently. By assigning higher values to these underexplored actions, UCB encourages the agent to investigate actions with greater uncertainty, characterized by larger confidence intervals. This exploration often results in a significant increase in cumulative rewards when these actions prove to be rewarding. However, as the agent gathers more data, the uncertainty about these actions decreases, and the confidence bounds shrink. Consequently, the agent shifts its focus away from exploration, leading to a decline in cumulative rewards when the agent learns that some actions were not as promising as initially expected. This process of early exploration followed by a correction phase accounts for the characteristic spikes in UCB's performance at the beginning of the learning process.