# CS 440/ECE 448 MP1: Part II

Sean Yen

# 1. WorldWithoutThief

## (a)

Following are the rewards for the first 10 episodes.

| | |
|---|---|
| 1 | −7.0 |
| 2 | −0.5 |
| 3 | −0.5 |
| 4 | 0.0 |
| 5 | 0.0 |
| 6 | 0.0 |
| 7 | 0.0 |
| 8 | 0.0 |
| 9 | 0.0 |
| 10 | 0.0 |

The reward remains at 0 for the rest of the simulation. This seems to be because the agent is "afraid" of trying new actions for fear of negative rewards. Here, we see that the agent starts by receiving a negative rewards, reduces that negative reward by a bit, and when it hits 0.0, it seems to settle for this being the best that it can do.

## (b)

Following are the rewards for the first 30 episodes.

| | |
|---|---|
| 1 | −21.5 |
| 2 | −10.0 |
| 3 | −11.5 |
| 4 | −7.0 |
| 5 | −13.0 |
| 6 | −8.0 |
| 7 | −6.5 |
| 8 | −4.5 |
| 9 | −5.5 |
| 10 | −9.0 |
| 11 | −9.0 |
| 12 | −3.0 |
| 13 | −5.0 |
| 14 | 0.0 |
| 15 | −5.5 |
| 16 | −5.0 |
| 17 | −0.5 |
| 18 | 2.0 |
| 19 | −7.5 |
| 20 | −5.0 |
| 21 | 2.0 |
| 22 | −4.0 |
| 23 | 3.5 |
| 24 | −3.0 |
| 25 | −4.5 |
| 26 | −3.5 |
| 27 | −6.0 |
| 28 | 2.5 |
| 29 | 1.0 |
| 30 | 17.0 |

We clearly see that the reward passes 0.0 and becomes positive in this case. In fact, in episode 31, we see that the reward jumps to 147.5! The reason that we are able to achieve such good rewards is because with $\epsilon = 0.1$, our agent is randomly selecting actions that may not be in the policy. This is allowing it to explore the "world" more fully to find more optimal solutions to the problems.

**(c)**

Following are the rewards for the last 10 episodes.

| | |
|---|---|
| 1 | −34.5 |
| 2 | −28.5 |
| 3 | −15.0 |
| 4 | −29.5 |
| 5 | −9.5 |
| 6 | −33.5 |
| 7 | −18.5 |
| 8 | −41.0 |
| 9 | −49.0 |
| 10 | −16.0 |

In this case, notice that even after 1000 episodes, the reward is still negative (and *very negative*). It seems that $\epsilon$ is set too high and the agent is taking random actions too often. This is causing it to choose random actions more often than it should and it is "stumbling" into a non-optimal path.

# 2. WorldWithThief

**(a)**

Following are the rewards for the first 10 episodes.

| | |
|---|---|
| 1 | −26.5 |
| 2 | −20.0 |
| 3 | −14.5 |
| 4 | −20.0 |
| 5 | −15.0 |
| 6 | −11.5 |
| 7 | −14.5 |
| 8 | −9.5 |
| 9 | −11.0 |
| 10 | −18.0 |

The rest of the episode rewards continue much like this. The agent rarely achieves positive rewards. It seems that without knowing where the thief is, since the thief moves randomly, when the agent encounters a thief, it can't figure out why it's incurring a negative reward and so it doesn't know how to avoid the thief effectively.

**(b)**

Following are the rewards for the first 10 episodes.

| | |
|---|---|
| 1 | −79.5 |
| 2 | −50.5 |
| 3 | −43.5 |
| 4 | −35.5 |
| 5 | −15.0 |
| 6 | 85.5 |
| 7 | 148.5 |
| 8 | 149.0 |
| 9 | 135.5 |
| 10 | 133.0 |

Over the course of the 1000 episodes, the agent goes on to achieve awards of around 300! In this case, since the agent knows where the thief is and has states to account for the position of the thief, it can find a deterministic way of avoiding the thief and delivering the goods to incur the greatest reward.

## (c)

See the following table of a survey of possible $\alpha \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ and $\epsilon \in \{0.0, 0.05, 0.1, 0.5, 0.7\}$ values and the resulting rewards at episode 1000.
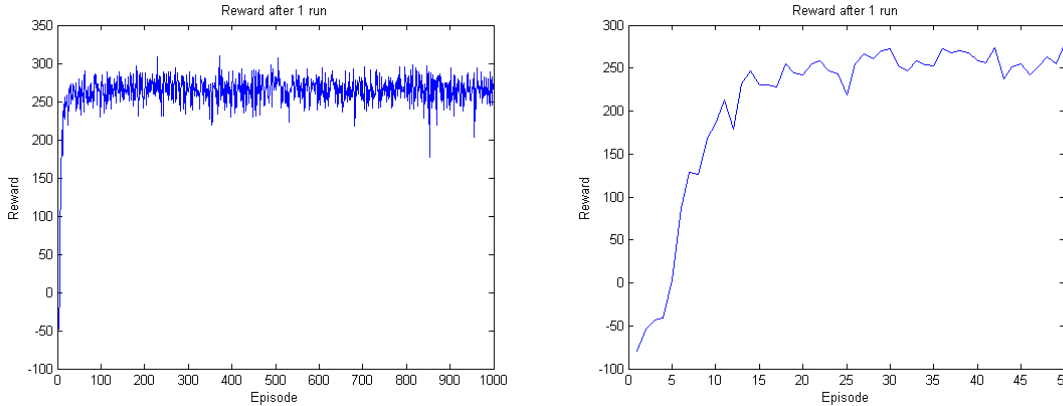
|  |  | $\epsilon$ | | | | |
|---|---|---|---|---|---|---|
|  |  | 0.0 | 0.05 | 0.1 | 0.5 | 0.7 |
|  | 0.1 | 77.0 | 267.5 | 200.5 | -192.5 | -213.5 |
|  | 0.3 | 111.0 | 281.5 | 219.5 | -222.0 | -247.0 |
| $\alpha$ | 0.5 | 109.0 | 286.0 | 186.5 | -212.5 | -270.5 |
|  | 0.7 | 107.0 | 222.5 | 176.0 | -187.5 | -269.5 |
|  | 0.9 | 76.0 | 178.5 | 142.5 | -200.5 | -250.0 |

The values that I chose cover the main support of $\alpha$ and $\epsilon$ and show what I think is a representative survey of the effect of these parameters on the reward at episode 1000.

Looking at these results, it's clear that choose $\epsilon$ as 0.5 or 0.7 clearly is not optimal. We get very large negative rewards in these cases. It seems like choosing a small, non-zero $\epsilon$ (0.05 seems to be the best) produces the best results. For $\alpha$, we see a more hyperbolic relationship. Small and large $\alpha$'s given worse results. This is likely because in these cases, we are trusting the new information too much (or not enough). The best $\alpha$'s appear to be 0.3 and 0.5.
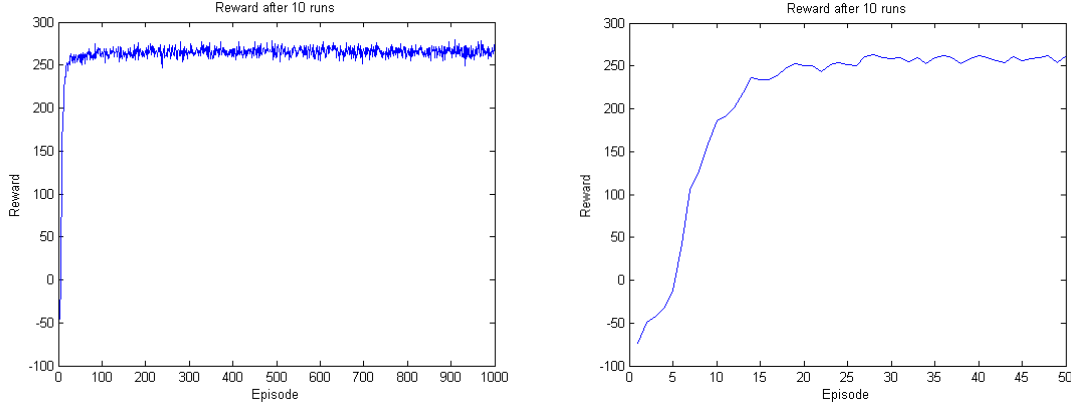
# Problem 3

The following plots were generated using $\alpha = 0.3$ and $\epsilon = 0.05$. These plots are episode by reward which shows the change in the reward as our agent refines its policy. The plot on the left will show rewards as they change over the course of 1000 episodes. The plot on the right shows rewards over the first 50 episodes where most of the growth will occur.



Notice that the reward seems to become level around 250 ($\mu_{\text{episode} \in [20:1000]} \approx 265$). In fact, within the first 20 episodes, the reward reaches this asymptote. Afterwards, the values hover around this asymptotic value ($\sigma_{\text{episode} \in [20:1000]} \approx 15$).

Next we consider the same plots but averaged over 10 runs of 1000 episodes.



We now observe the same general behavior, but the function seems to be much smoother. This can also be seen in the parameters $\mu_{\text{episode}\in[20:1000]}$ and $\sigma_{\text{episode}\in[20:1000]}$. In this case, $\mu_{\text{episode}\in[20:1000]} \approx 265$ and $\sigma_{\text{episode}\in[20:1000]} \approx 5$. It's notable that the standard deviation, $\sigma$, is much lower than it was in the previous case when we just considered 1 run of 1000 episodes. This is due to the fact that averaging here reduces the noise in the estimate of our behavior.