

# Project 2

Jerry Duncan

November 5, 2020

## 1 Introduction

For this project, we are taking a look at a simple Reinforcement Learning problem where there exists a small,  $dxd$  grid and there exists a robot and bomb somewhere on this grid. Our goal is to teach the robot how to get the bomb off the grid as quickly as possible. We will apply two different methods to this problem, Q-Learning and First-Visit Monte Carlo, to see which performs better at a) generally solving the problem and b) which does so in the fastest manner. First we will form the problem as an MDP, where the state consists of the robot and bomb's position. Then we will apply those two methods using two different reward structures and 2 different behavior policies in the case of Q-Learning, to see how well they solve the problem. And lastly, we will compare the time it takes for them to converge to their solutions.

**Important Note** I set up my problem so that if we don't push the bomb out by step 1,000, we lose the problem. This means that we will still use the episode to learn for the MC Learning because this greatly simplified the programming for the problem. It still converges and works, but it does mean that some of the Q values are a lot larger than might be expected when talking about them in the Results.

## 2 Markov Decision Process

**State space** The state space for this problem is very simple. The environment must keep track of two things and allow the robot to know of those

two things: robot X, robot Y, bomb X, and bomb Y. This means the dimensionality of the state space ends up being  $d^4$  and since  $d = 8$  then the state space has 4096 total states.

**Action space** The only actions our robot can take at any state are the cardinal directions  $N, E, S, W$ . So our action space  $A = N, E, S, W$  for each and every state.

**Reward Structure** We want to test out two different reward structures, both of which are given to us. The first is to always give a -1 reward regardless of move. The other, is more complicated, but gives +10 when the robot puts the bomb in the river, +1 when the robot moves the bomb further away from the center of the grid, and -1 for any other move.

### 3 Code and Data Structures

I used a class-based structure for this project to keep all data close to its source. The environment is modeled as a class that has the following methods:

- *start*: Initializes a random starting state to begin an episode
- *move*: Takes an action, given the current state, and moves the robot (and bomb if necessary), as well as returning the reward for that action
- *reward*: Determines the reward based on a given state and action and which reward structure being used
- Helpers
  - *dist*: Calculates the distance from the center of the board to an arbitrary position
  - *next\_position*: Calculates the next position given a current position and an action

The `MC_Learner` and `Q_Learner` classes subclass `Learner`. All Learners are structured as follows:

- *\_\_init\_\_*: Initializes an array for Q values and for the policy

- *learn*: Generates episodes and learns on them, depending on the type of learning (Q or MC)
- *run*: Runs the learning function  $n$  times.
- *plot*: Takes a starting state, generates an episode, and plots it
- Helpers
  - *final\_policy*: Prepares the learned policy for use in episodes after training
  - *q-values*: Runs an episode and prints out the Q values

## 4 Results

All experiments are performed over 50,000 episodes.

Before we start looking at results and discussing why they make sense, let's first determine what should the optimal route given the constant or dynamic reward structure be?

Given the board from the assignment, we as humans can see the optimal path has to be 5 moves long because there is no shorter path that can get the bomb off the grid. And for the constant reward structure, the optimal path overall is the same as the optimal path for maximizing the reward. Specifically, that sequence of actions would be: East, East, North, West, and West for a total of -5 reward. However, for the dynamic reward structure, we get +1 when we move the bomb further away from the center. This means that we can potentially make up for some of our earlier -1 moves by taking some +1 moves. While it may not be obvious at first glance, I propose that the following sequence of actions would product a higher reward than the ones used for the constant reward structure: East, North, North, North, North, North. This set of actions gives us a reward of 11 versus 8 for the action set from the optimal constant reward path. This is because while the first two actions are normal steps and we receive -2, the third, fourth, and fifth action all move the box further away from the center giving us +3 and the 6th final action gives us +10 for removing it from the board for a total of 11. It also contains 6 actions, meaning it's impossible to be the optimal for the constant reward structure. I propose that since these two paths are the optimal, they should be what we expect our models to learn. Make note

that the 6-move optimal path for maximizing the reward for the dynamic reward structure is not optimal for the problem itself, but this is a failing of the dynamic reward structure.

#### **Part a Q-Learning with a Constant Reward and a Uniform Policy**

- $Q((\text{robotX: 0, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -7.00, E: -5.00, S: -7.00, W: -6.00}]$
- $Q((\text{robotX: 1, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -5.00, E: -4.00, S: -6.00, W: -6.00}]$
- $Q((\text{robotX: 2, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -3.00, E: -5.00, S: -5.00, W: -5.00}]$
- $Q((\text{robotX: 2, robotY: 4}), (\text{bombX: 1, bombY: 4})) = [\text{N: -4.00, E: -4.00, S: -4.00, W: -2.00}]$
- $Q((\text{robotX: 1, robotY: 4}), (\text{bombX: 0, bombY: 4})) = [\text{N: -3.00, E: -3.00, S: -3.00, W: -1.00}]$

We end up with nice, round whole numbers because our policy has forced us to take each action equally, meaning that our estimates of the Q's are more likely to be accurate.

#### **Q-Learning with a Constant Reward and a Greedy Policy**

- $Q((\text{robotX: 0, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -4.90, E: -4.81, S: -4.94, W: -4.85}]$
- $Q((\text{robotX: 1, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -4.08, E: -4.00, S: -4.28, W: -4.18}]$
- $Q((\text{robotX: 2, robotY: 5}), (\text{bombX: 1, bombY: 4})) = [\text{N: -3.00, E: -4.01, S: -4.54, W: -3.43}]$
- $Q((\text{robotX: 2, robotY: 4}), (\text{bombX: 1, bombY: 4})) = [\text{N: -3.99, E: -3.96, S: -3.98, W: -2.00}]$
- $Q((\text{robotX: 1, robotY: 4}), (\text{bombX: 0, bombY: 4})) = [\text{N: -3.00, E: -3.00, S: -3.00, W: -1.00}]$

While we still end up with the right path when using a greedy policy, the estimated Q values are less accurate because despite going through 50,000 episodes, the chances that we visit these specific states and don't pick the greedy action are still relatively low. If we had trained longer, these would have converged further, nevertheless they're still good enough to find the optimal solution for this reward structure.

#### **Monte Carlo First-Visit with a Constant Reward**

- $Q((\text{robotX: } 0, \text{robotY: } 5), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -198.30, \text{E: } -203.12, \text{S: } -12.68, \text{W: } -183.27]$
- $Q((\text{robotX: } 0, \text{robotY: } 6), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -95.39, \text{E: } -73.63, \text{S: } -10.66, \text{W: } -58.54]$
- $Q((\text{robotX: } 0, \text{robotY: } 7), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -163.37, \text{E: } -8.65, \text{S: } -136.27, \text{W: } -126.16]$
- $Q((\text{robotX: } 1, \text{robotY: } 7), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -257.84, \text{E: } -6.83, \text{S: } -165.50, \text{W: } -118.29]$
- $Q((\text{robotX: } 2, \text{robotY: } 7), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -5.59, \text{E: } -150.20, \text{S: } -210.05, \text{W: } -212.83]$
- $Q((\text{robotX: } 2, \text{robotY: } 6), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -5.01, \text{E: } -119.67, \text{S: } -182.41, \text{W: } -142.35]$
- $Q((\text{robotX: } 2, \text{robotY: } 5), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -3.35, \text{E: } -102.30, \text{S: } -64.37, \text{W: } -88.88]$
- $Q((\text{robotX: } 2, \text{robotY: } 4), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -33.09, \text{E: } -54.13, \text{S: } -39.45, \text{W: } -2.06]$
- $Q((\text{robotX: } 1, \text{robotY: } 4), (\text{bombX: } 0, \text{bombY: } 4)) = [\text{N: } -4.34, \text{E: } -3.61, \text{S: } -6.67, \text{W: } -1.00]$

Because the episodes are so long, our initial estimates of Q can be very large. Especially because each state is so rare to occur and exploration is still somewhat rare.



Figure 1: The path taken by the Q-Learning robot with a Constant Reward and a Uniform Policy. Blue-green is the robot. Yellow is the bomb.

**Part b** The path taken by **Q-Learning with a Constant Reward and a Uniform Policy** is shown in Figure 4. It makes sense that we found the optimal path because our Q estimates have converged perfectly for this scenario and our final policy can easily find this path. The path taken by **Q-Learning with a Constant Reward and a Greedy Policy** is shown in Figure 4. Similar to the uniform policy Q learner, this one has near-perfect Q estimates and can still easily find the path. And the path taken by **Monte Carlo First-Visit with a Constant Reward** is shown in Figure 4. The Monte Carlo path is a bit rough and nowhere near optimal. The reason for this is likely because the first time we solved this scenario in our training, it was nonoptimal and we didn't explore enough to find better Q estimates.

**Part c** **Q-Learning with a Dynamic Reward and a Uniform Policy**

- $Q((\text{robotX: } 0, \text{ robotY: } 5), (\text{bombX: } 1, \text{ bombY: } 4)) = [N: 10.00, E:$



Figure 2: The path taken by the Q-Learning robot with a Constant Reward and a Greedy Policy. Blue-green is the robot. Yellow is the bomb.

MC 50k First-Visit Constant Reward

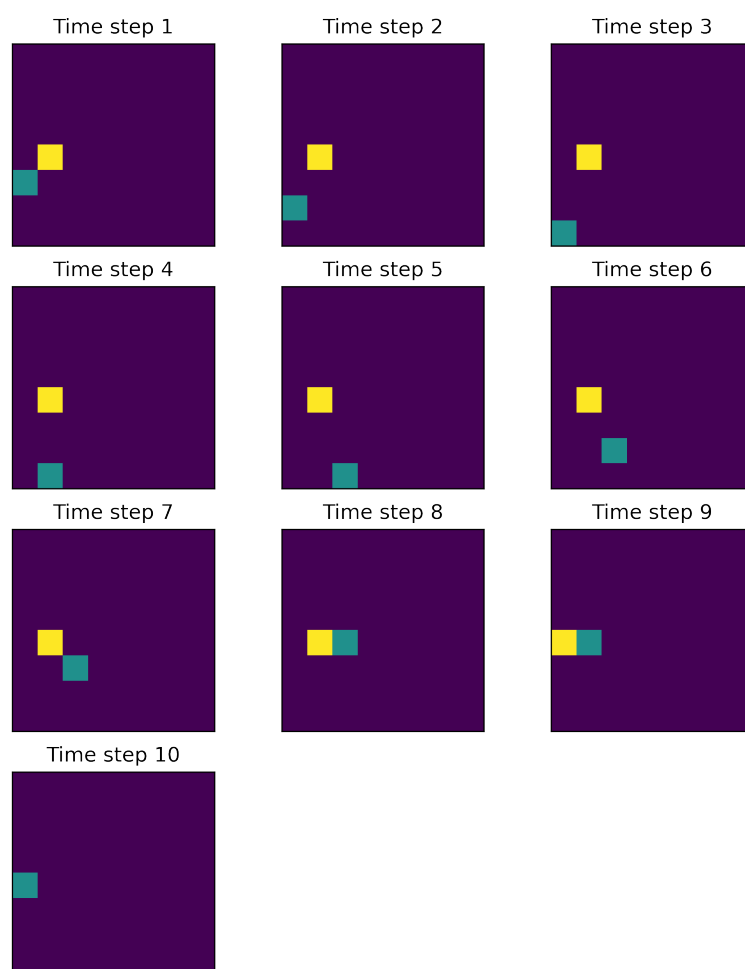


Figure 3: The path taken by the First-Visit Monte Carlo robot with a Constant Reward. Blue-green is the robot. Yellow is the bomb.



11.00, S: 9.00, W: 10.00]

- $Q((\text{robotX: } 1, \text{ robotY: } 5), (\text{bombX: } 1, \text{ bombY: } 4)) = [\text{N: } 12.00, \text{ E: } 10.00, \text{ S: } 10.00, \text{ W: } 10.00]$
- $Q((\text{robotX: } 1, \text{ robotY: } 4), (\text{bombX: } 1, \text{ bombY: } 3)) = [\text{N: } 13.00, \text{ E: } 11.00, \text{ S: } 11.00, \text{ W: } 11.00]$
- $Q((\text{robotX: } 1, \text{ robotY: } 3), (\text{bombX: } 1, \text{ bombY: } 2)) = [\text{N: } 12.00, \text{ E: } 10.00, \text{ S: } 10.00, \text{ W: } 10.00]$
- $Q((\text{robotX: } 1, \text{ robotY: } 2), (\text{bombX: } 1, \text{ bombY: } 1)) = [\text{N: } 11.00, \text{ E: } 9.00, \text{ S: } 9.00, \text{ W: } 9.00]$
- $Q((\text{robotX: } 1, \text{ robotY: } 1), (\text{bombX: } 1, \text{ bombY: } 0)) = [\text{N: } 10.00, \text{ E: } 9.00, \text{ S: } 8.00, \text{ W: } 8.00]$

Similar to the constant reward Q's, because the uniform policy causes us to explore constantly, our Q estimates become perfect. In a way, because this has run so long, the forced exploration has essentially solved this problem for all states.

### **Q-Learning with a Dynamic Reward and a Greedy Policy**

- $Q((\text{robotX: } 0, \text{ robotY: } 5), (\text{bombX: } 1, \text{ bombY: } 4)) = [\text{N: } 0.12, \text{ E: } 4.71, \text{ S: } -0.20, \text{ W: } 0.40]$
- $Q((\text{robotX: } 1, \text{ robotY: } 5), (\text{bombX: } 1, \text{ bombY: } 4)) = [\text{N: } 4.66, \text{ E: } 9.00, \text{ S: } -0.09, \text{ W: } 0.16]$
- $Q((\text{robotX: } 2, \text{ robotY: } 5), (\text{bombX: } 1, \text{ bombY: } 4)) = [\text{N: } 10.00, \text{ E: } 4.23, \text{ S: } 4.03, \text{ W: } 3.22]$
- $Q((\text{robotX: } 2, \text{ robotY: } 4), (\text{bombX: } 1, \text{ bombY: } 4)) = [\text{N: } 8.94, \text{ E: } 8.84, \text{ S: } 8.95, \text{ W: } 11.00]$
- $Q((\text{robotX: } 1, \text{ robotY: } 4), (\text{bombX: } 0, \text{ bombY: } 4)) = [\text{N: } 7.99, \text{ E: } 7.94, \text{ S: } 7.95, \text{ W: } 10.00]$

Similar to the greedy policy from before, the Q's are not nearly perfect estimates, however they highest Q's do correspond to the correct actions. While we explore a little bit, we're able to get a "good-enough" estimate for Q very quickly. And once we get that "good-enough" estimate, we don't

explore enough to find better estimates for the other actions. This tends to be okay though because all that matters is that our estimate for the optimal action is the best, which often happens. Do note that we didn't find the optimal path for the dynamic reward though, likely due to not exploring enough.

### Monte Carlo First-Visit with a Dynamic Reward

- $Q((\text{robotX: } 0, \text{robotY: } 5), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } -127.28, \text{E: } 6.10, \text{S: } -78.30, \text{W: } -81.98]$
- $Q((\text{robotX: } 1, \text{robotY: } 5), (\text{bombX: } 1, \text{bombY: } 4)) = [\text{N: } 6.82, \text{E: } -18.75, \text{S: } -29.49, \text{W: } -32.04]$
- $Q((\text{robotX: } 1, \text{robotY: } 4), (\text{bombX: } 1, \text{bombY: } 3)) = [\text{N: } 2.53, \text{E: } 7.84, \text{S: } 0.09, \text{W: } 4.75]$
- $Q((\text{robotX: } 2, \text{robotY: } 4), (\text{bombX: } 1, \text{bombY: } 3)) = [\text{N: } 9.31, \text{E: } 1.47, \text{S: } 4.97, \text{W: } 7.26]$
- $Q((\text{robotX: } 2, \text{robotY: } 3), (\text{bombX: } 1, \text{bombY: } 3)) = [\text{N: } 5.91, \text{E: } 6.12, \text{S: } 7.30, \text{W: } 10.61]$
- $Q((\text{robotX: } 1, \text{robotY: } 3), (\text{bombX: } 0, \text{bombY: } 3)) = [\text{N: } 7.88, \text{E: } 7.54, \text{S: } 7.51, \text{W: } 10.00]$
- $Q((\text{robotX: } 2, \text{robotY: } 3), (\text{bombX: } 0, \text{bombY: } 3)) = [\text{N: } -34.73, \text{E: } -33.75, \text{S: } -93.51, \text{W: } 8.79]$
- $Q((\text{robotX: } 1, \text{robotY: } 3), (\text{bombX: } 0, \text{bombY: } 3)) = [\text{N: } 7.88, \text{E: } 7.54, \text{S: } 7.51, \text{W: } 10.00]$

Similar to the MC learner with the constant reward, our estimates are very far off. Again this can be attributed to the lack of discounting factor, the possibility of staying with the first episode that terminates from that position, and the lack of exploration.

Overall, the reward structure has very little to do with the final Q estimates other than their values (being positive / negative). In all three cases, each method learned roughly the same across the two reward structures.

The path taken by **Q-Learning with a Dynamic Reward and a Uniform Policy** is shown in Figure 4. It makes sense that we found the optimal path because our Q estimates have converged perfectly for this

scenario and our final policy can easily find the optimal 6-move path for the dynamic reward scenario. The path taken by **Q-Learning with a Dynamic Reward and a Greedy Policy** is shown in Figure 5. Similar to the uniform policy Q learner, this one has near-perfect Q estimates and can still easily find a near-optimal path, but it is likely not finding the optimal path due to not exploring enough. And the path taken by **Monte Carlo First-Visit with a Dynamic Reward** is shown in Figure 6. The Monte Carlo path closer to a near optimal path this time. Again, the reason for this is partially due to a lack of exploration and partially due to the chances of finding a short episode (and therefore smaller Q values) is very unlikely.

Overall, forcing the learners to explore, while slow, does seem to produce optimal solutions.

**Part d** Because we went through 50k episodes, the graphs for rewards per episode is messy. A better graph would be the average reward over runs.

In Figure 7, I show the reward per episode using a constant reward structure for all three methods: MC Learning, Q-Learning with a Greedy Policy, and Q-Learning with a Uniform Policy. As we can see, it's fairly messy to understand. I've devised a slightly different graph to better show the reward per episode that I deem the average reward per run which is a rolling average calculated like so, where  $i$  is the episode number:

$$Avg(i) = \frac{\sum_1^i R_i}{i}$$

In Figure 8, I show the average reward per episode using a constant reward structure. From these two graphs, we can see that the episodes generated Q-Learning with a Greedy Policy have a very high reward. Monte Carlo achieves similar results, but it takes a lot more episodes to get there. It's harder to determine the learning of Q-Learning with a Uniform Policy because the episodes it generates are always random.

In Figure 9, I show the reward per episode using a dynamic reward structure for all three methods: MC Learning, Q-Learning with a Greedy Policy, and Q-Learning with a Uniform Policy. Similar to the constant reward structure results, it's also messy to understand. So in Figure 10, I show the average reward per episode using a dynamic reward structure. From these two graphs, we can see that the episodes generated Q-Learning with a Greedy Policy have a very high reward and learn very very fast compared to

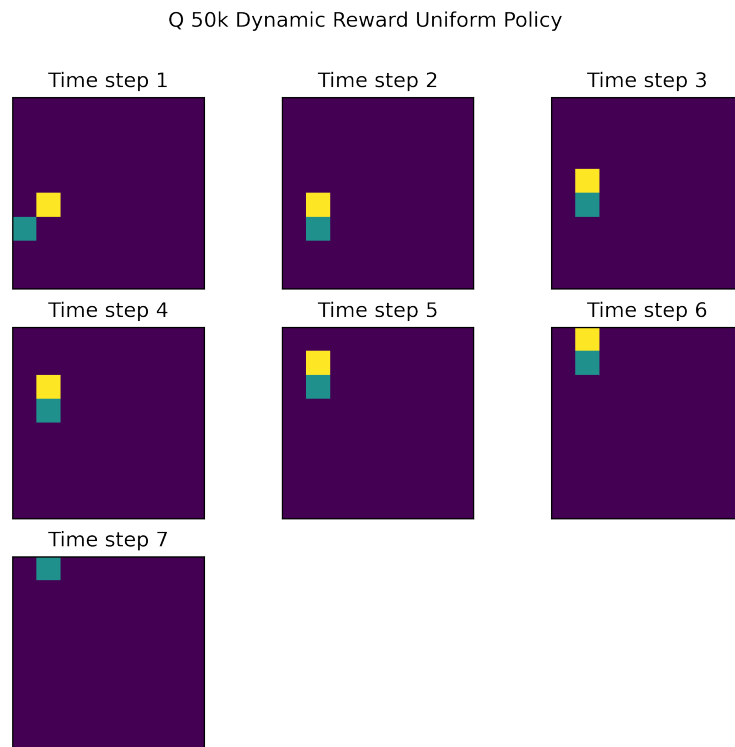


Figure 4: The path taken by the Q-Learning robot with a Dynamic Reward and a Uniform Policy. Blue-green is the robot. Yellow is the bomb.



Figure 5: The path taken by the Q-Learning robot with a Dynamic Reward and a Greedy Policy. Blue-green is the robot. Yellow is the bomb.

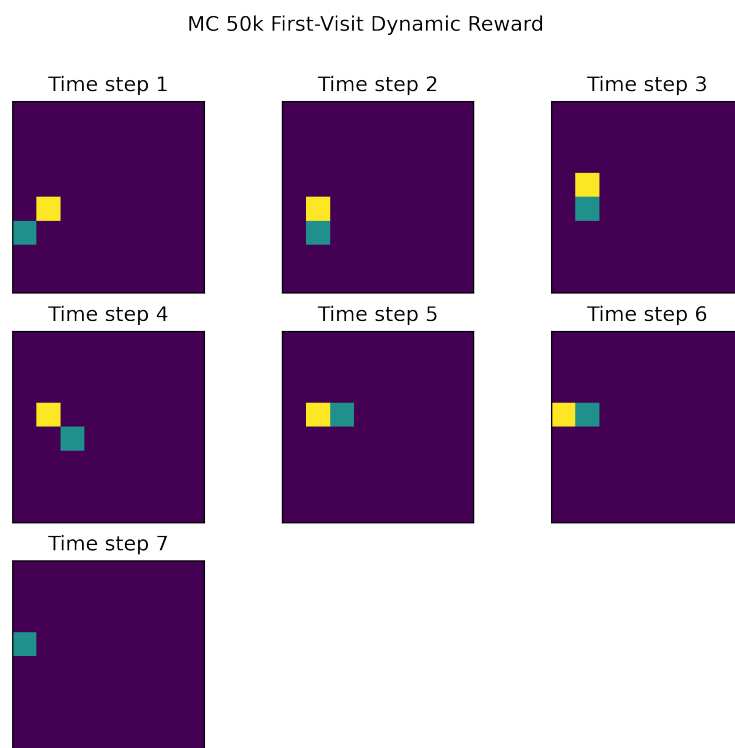


Figure 6: The path taken by the First-Visit Monte Carlo robot with a Dynamic Reward. Blue-green is the robot. Yellow is the bomb.

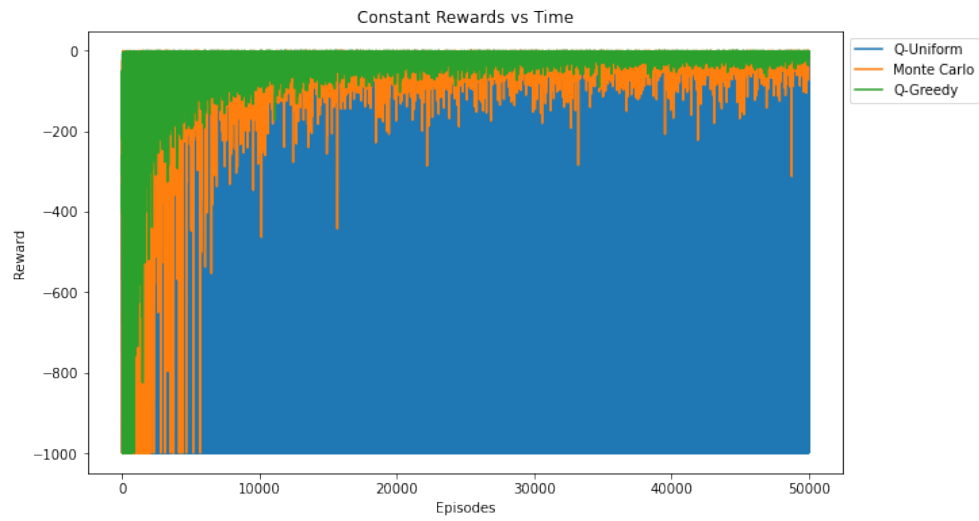


Figure 7: Rewards per episode using a constant reward structure.

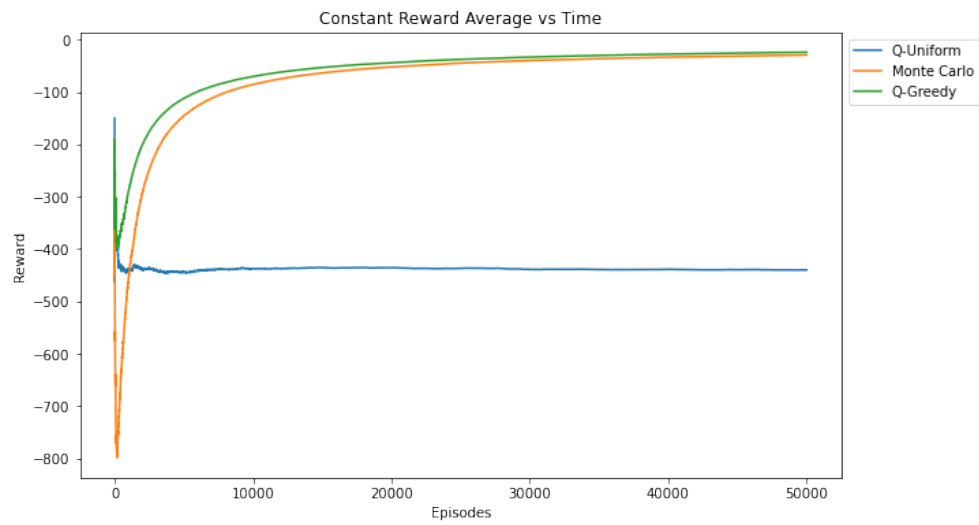


Figure 8: Average reward per episode using a constant reward structure.

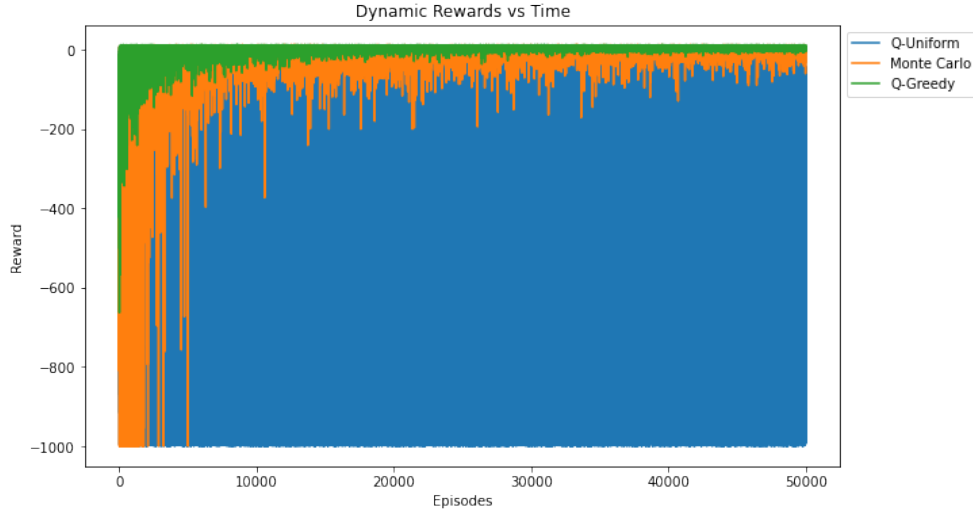


Figure 9: Rewards per episode using a dynamic reward structure.

when using a constant reward structure. Monte Carlo achieves similar results when using either reward structure. And again it's still harder to determine the learning of Q-Learning with a Uniform Policy because the episodes it generates are always random.

Both of these sets of results make sense because regardless of reward structure, Monte Carlo takes a good amount of successful episodes to learn something useful and that's why the results for Monte Carlo are roughly the same across both. For Q-Learning, the dynamic reward structure gives us better estimates short term for the cases where we move the bomb further away from the center and prevent Q values from exploding when we have an unsuccessful episode.

**Part e** In Figure 11 I show the time taken up to each episode versus the number of episodes ran using a constant reward. In Figure 12 I show the same graph but for dynamic rewards. In both cases, the Q-Learner with the uniform policy takes the longest — this makes sense because each episode is essentially random, we aren't using any learned knowledge to dictate our episodes so learning takes a long time. In the case of Monte Carlo and Q-Learning with the greedy policy start to go a lot faster once they learn a near-optimal policy because they're using their learned knowledge to generate episodes.



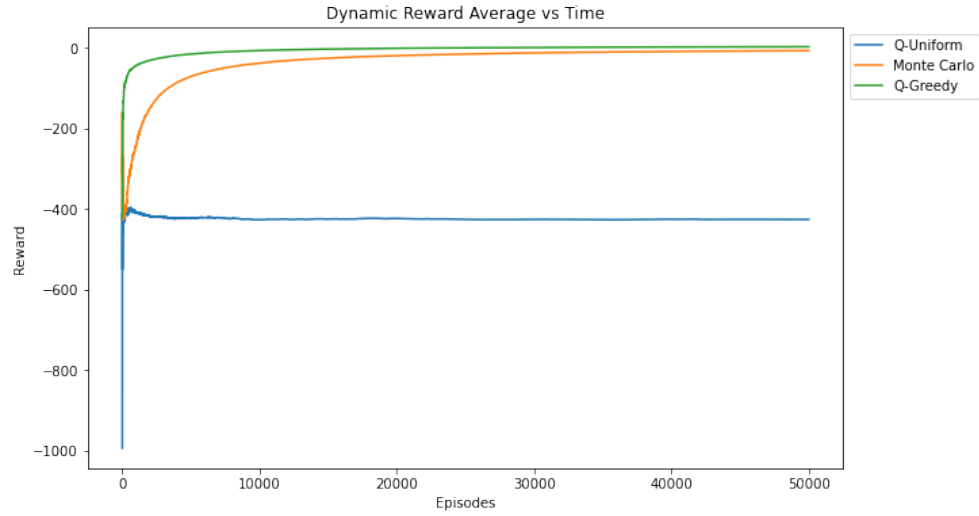


Figure 10: Average reward per episode using a dynamic reward structure.

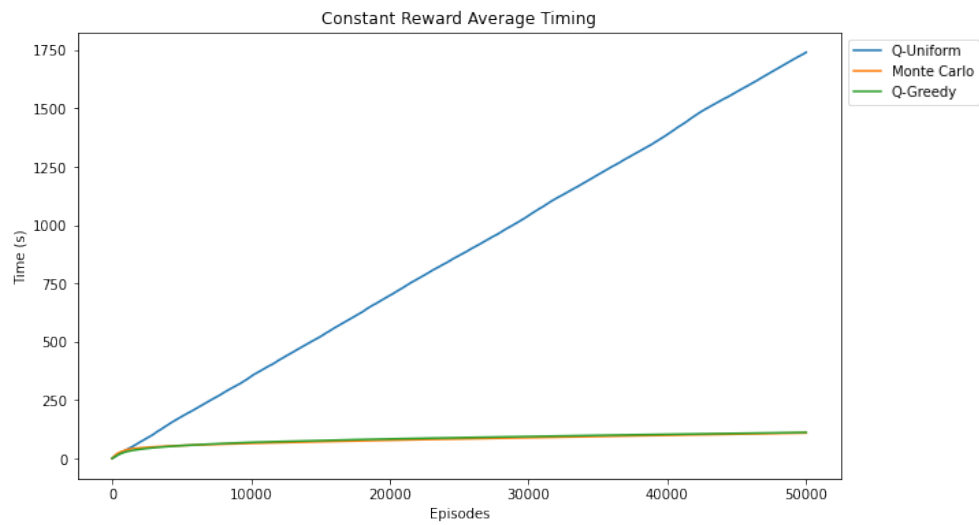


Figure 11: Total time taken for each learning method vs episodes using a constant reward structure.

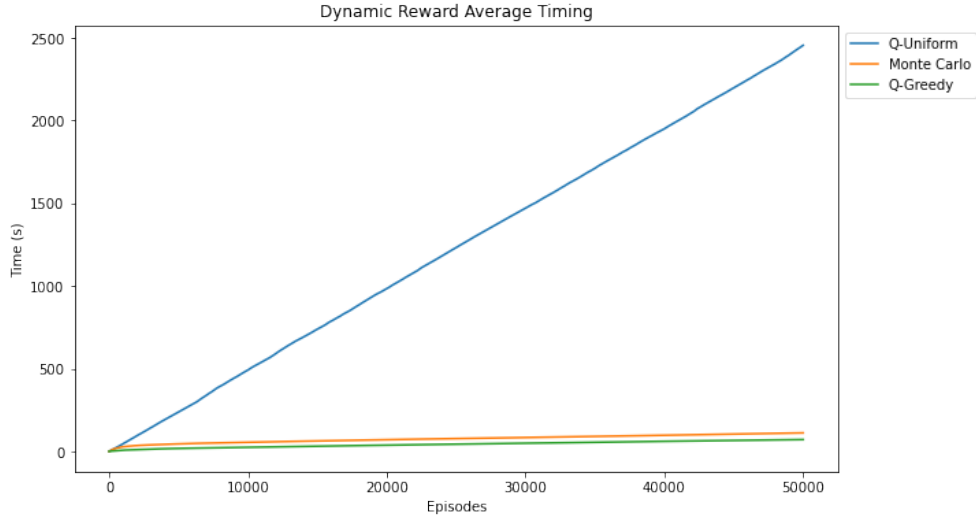


Figure 12: Total time taken for each learning method vs episodes using a dynamic reward structure.

## 5 Conclusion

Overall we've tested 3 different learning methods using 2 different reward structures to see how well they perform in both finding the optimal solutions to the robot-bomb problem and how long they take to train. From these results we can conclude that MC learning, while able to find a near-optimal policy, performs worse than both Q-learners in performance, and only better than Q-Uniform in terms of time taken. And between the two Q-learners, the one with the uniform policy ends up with the better final policy than the one using the greedy policy, but it takes much longer to train. The main take away is that each method has its pros and cons and are all viable candidates at solving basic Reinforcement Learning problems.