

---

# Final Report: Malmo Minecraft Mazerunner

---

## Authors

Jordan Choi  
Yixuan Li  
Joshua Lu

## Abstract

Our project is to develop and train an agent that can guide our Minecraft user from the beginning of a maze to the end. The main idea is to create an agent that is capable of learning and familiarizing itself in different environments and then make the most optimal decision at every move to find the best overall combination of actions to solve the maze. The agent is first placed in a predetermined starting position and its goal is to reach the diamond block whose location is unknown to the agent and is up to it to find that ending point. The agent will decide between a combination of simple actions such as moving as well as complicated actions like jumping in order to reach the endpoint of the maze.

We have accomplished two of the three goals we had set out for our project. Firstly, our minimum goal was to accomplish simple movement from the start point to the end point of a basic maze (small in size and not difficult to navigate). Ever since the beginning, we have been using the deep Q-learning approach to simulate reinforcement learning and the agent showed capability in solving basic mazes efficiently by finding the fastest path from start to the end of the maze. Moving on, we completed our intermediate goal of incorporating more complicated actions into our learning agent; the ability to jump. Our agent was able to use a variety of different combinations of the now simple and complicated actions to solve the mazes we fed it. We had originally planned for our hard goal to be the addition of further complicated actions which would include digging and building, but due to how our progress turned out for the previous two goals, we decided that it would be best to focus honing our agent to complete our minimum and intermediate goals with no flaws.

## 1 Introduction

The goal of our project is to train our agent to be able to navigate a maze from beginning to end. We want our agent to be able to learn to direct the user through varying obstacles based on its experience amassed from reinforced learning. What separates an artificial intelligence agent from just a well-coded maze solver is its ability to learn and adapt and make good decisions in unknown environments. We wanted to create an agent that isn't just good at solving mazes through countless repeated runs because any computer program could do that. We wanted an agent that can make smart decisions at every move and find the optimal route to navigate itself through an unfamiliar maze without running it an unlimited number of times. We decided to implement a deep Q-learning approach because it addresses both of these demands. Q-learning is interesting because it implements randomization during the testing phase to force the algorithm to deal with unpredictable states and require a degree of generalization when it comes to decision making. This answers our problem as implementing Q-learning allows our agent to rely on reinforced learning to be able to build up enough experience to make smart decisions when placed in a foreign maze. We played around with different approaches of using grid-based observations using tools provided by the Malmo platform

39 versus frame-based observations where each frame of the game represented a different state and is  
40 visual dependent as compared to the frame-based approach. At the moment, we have completed our  
41 minimum goal using both approaches. Continuing on to our second milestone, we will continue  
42 to analyze trade offs and make a decision on which approach to stick with as we continue to try to  
43 complete more complex objectives.

## 44 2 Background

45 Our project is set in Minecraft, a video game that allows a player to control an agent to explore and  
46 craft in the world, and build complicated structures. We are using Project Malmö, an open source  
47 platform built on Minecraft for artificial intelligence experimentation to help the agent act and sense  
48 within the Minecraft environment. We are also using its provided Q-learning library.

## 49 3 Problem Statement

50 The agent first starts at a predetermined starting block in a maze that also has a goal block. The maze  
51 has redstone blocks for “walls” and quartz blocks as possible paths for the agent to take. The agent  
52 has no information on the environment at the start but can slowly figure out its environment as it  
53 explores to find and reach the goal state.

54 The agent was given five actions: the agent can move forwards, backwards, turn left, turn right, or  
55 jump in a direction. We used Malmö to control the agent’s movement and actions. Furthermore, we  
56 assigned the algorithm this reward system:

- 57 • -1 for each move in a direction
- 58 • -1 for each jump
- 59 • -100 for dying
- 60 • +500 for getting to the gold block
- 61 • +1000 for getting to the diamond block
- 62 •  $-10 \times (\text{distance from goal})$  for each episode

63 The minimum goal of our project was to accomplish simple movement from the start point to the end  
64 point in a basic maze. The first milestone was using only back and forth movement and then in the  
65 second milestone, use a combination of front, back, left, and right movement to navigate to the end  
66 point. The realistic goal will be to add in jumping/digging actions in order for the agent to achieve its  
67 goal. The first milestone will be to teach the agent to learn to dig when able. The second milestone  
68 was to teach the agent to jump when it is more efficient to do so, rather than exploring and using a  
69 longer path. The ambitious goal will be to allow the agent to build in the maze to help it reach the  
70 end point. The agent should learn whether it would be more efficient to build and jump over a wall  
71 rather than going around the maze walls. For this goal, its milestone would be to build to jump over a  
72 wall or dig under the wall rather than going around the wall to get to the end point.

## 73 4 Method

74 Our agent follows the Markov Decision Process to determine the most efficient path. At each step,  
75 the agent examines and interacts with the environment and receives a certain reward. With each  
76 action, the agent receives a reward that depends on all previous states and actions it took to get to the  
77 current state. Thus, the agent chooses the appropriate action for its current state depending on actions  
78 in previous states.

79 Now given we know the expected reward of each action at every step, we use Q-learning, to generate  
80 the maximum total reward, which in Q-learning is called the Q-value.

$$Q(s, a) = r(s, a) + \max_a Q(s', a)$$

81 Figure 1: Markov Decision Process: Reinforcement Learning

82 The Q-value we calculate comes from being at the state “s” at which we will make the action “a”. To  
 83 get Q, we will calculate the immediate reward of making action “a” at state “s”, shown by “r(s,a).”  
 84 We will add this with the highest Q-value possible from the next state s’. We include the value of  
 85 gamma as it controls how much influence the future reward has by either increasing or decreasing it  
 86 based on what value we assert gamma as.  
 87 We also implement an  $\epsilon$ -greedy policy to avoid getting stuck in the local minimum if it only chooses  
 88 actions that optimize the Q-function. The  $\epsilon$ -greedy method will either choose the action with the  
 89 highest q-value or a random action. This decision is made randomly and is based on the value of  
 90 epsilon where the value of  $\epsilon$  is purposely set so that we lean towards using random actions at first,  
 91 but as training progresses, become more reliant on maximum Q-values instead. With this policy, the  
 92 agent will randomly choose an action that is not optimal so the agent has the chance to explore the  
 93 environment more at the off chance that the true optimal path is elsewhere and prevent overfitting.

## 94 5 Experiments

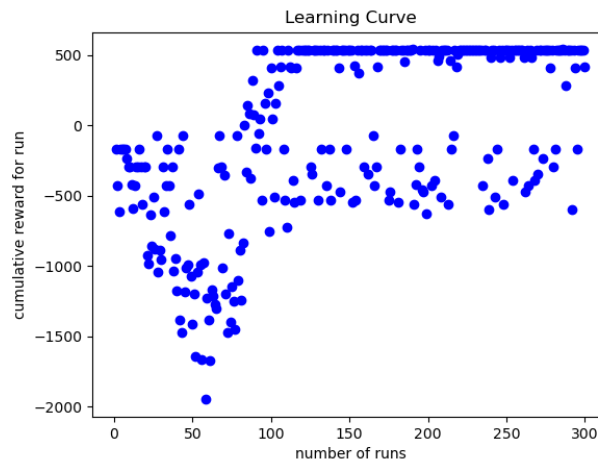
95 We first designed a map to see whether model would learn a simple map with a basic path get towards  
 96 the goal. This was to see whether the Q-learning algorithm was properly inserting value into the table.  
 97 To see the Q-learning process in action, we explicitly designed two maps that would demonstrate  
 98 whether the model would find the most optimal path. The first map had two goal nodes with two  
 99 simple paths with one path towards a gold block and one longer path towards a diamond block. The  
 100 second map had one goal node towards a diamond block with two paths. One that required jumping  
 101 that would ultimately take a shorter number of actions to get towards the goal and one that was much  
 102 longer that did not require jumping. Our final map was a more ambitious map that incorporated  
 103 multiple floors that required jumping, multiple goal nodes, and multiple paths. The goal of this map  
 104 was to really test the model, see how it would behave, and see if it would be able to find the most  
 105 optimal path.

## 106 6 Results

107 Results were measured as the cumulative reward for each episode over the number of episodes. An  
 108 expected result should show a positive relationship between the two variables. At first, the graph  
 109 should initially dip into the negative region as the agent is penalized for movements and deaths. As  
 110 the model continues exploring and exploiting, the model should eventually increase and peak at the  
 111 optimal solution. Once the model peaks, slight deviations should show a scatter of points under the  
 112 peak as a result of a positive non-zero epsilon value.

113 The first map reflects the expected outcome of the graph without any unforeseen trends. As the first  
 114 map is a simple path towards a goal, the agent finds the optimal path relatively early in training  
 115 showing an initial negative trend and then switching to a positive trend after discovering the goal  
 116 node for the first time.

117



118

Figure 2: Map 1 Plot

In the second map, the agent was quickly able to find the gold block early on as it was the very close to the starting node. This is reflected in the early trend line of consistent rewards of 220. The agent continues exploring and was able to find a more optimal path soon after discovering the gold block that the diamond block was more optimal and the trend quickly shifted, favoring a higher reward. This is reflected in the greater concentration of plots in the 800 range.

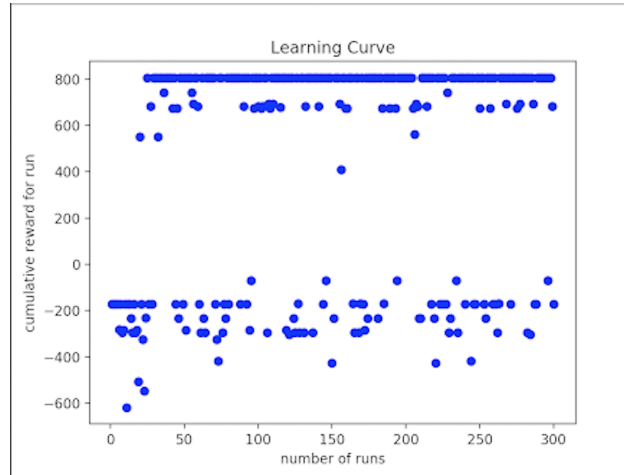


Figure 3: Map 2 Plot

A similar trend from the second map, follows in the third map albeit slightly less obvious at first. At first, the agent learns how to reach the goal node without jumping however, it later discovers that jumping over a block would produce a more optimal, shorter path due to the cost associated with each step taken. The graph reflects this trend rather subtly between episodes 60 and 120 versus episodes 120 to the end. The first half shows a slightly smaller reward of roughly 7 steps compared to the slightly greater second half. This is a direct result of the number of steps the agent took to move around a hazard versus saving those steps by jumping over the hazard. Exploration of the map resulted in expected trends both before the agent began consistently discovering a path to the goal node and after the agent discovered a path to the goal node

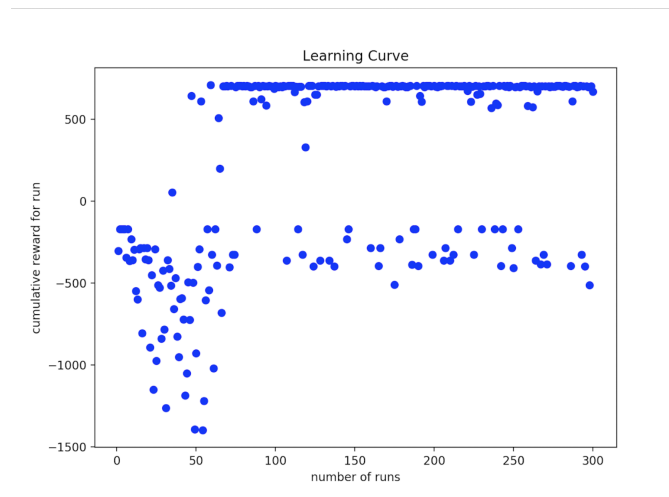


Figure 4: Map 3 Plot

The final map we tested involved a heavily more ambitious path to reach the optimal goal node. It was slightly larger and had multiple floors with multiple goal states. This map was created with the

142 intention of placing the agent under more realistic conditions and stresses in the world of Minecraft.  
 143 The resulting graph had the expected trends of exploration plots due to the epsilon being greater  
 144 than zero however, the two concurrent trend lines of the graph were difficult to interpret. The lines  
 145 represented the ending states of the agent learning how to reach the lesser reward of the gold block  
 146 however, it pointed out different paths of reaching the gold block as well as pointing out a bug within  
 147 our code.  
 148 In Malmo, as soon as the agent touches the block with any part of its body, the goal is considered  
 149 reached and the episode terminates. In the -485 cumulative reward range, the agent had learned  
 150 to reach the goal node by climbing stairs and walking to the goal node. This state was further by  
 151 roughly 16 blocks from the most optimal diamond block goal of which we were using to calculate  
 152 the `distance_from_goal` value. In the 10 cumulative reward range, the agent had learned to reach  
 153 a block one floor beneath the goal we had intended for and jumped to touch its head to the goal; an  
 154 unforeseen solution to reach the goal. This state was further by roughly 10 blocks from the most  
 155 optimal goal node. This discrepancy of roughly 6 blocks explains the two trends when calculating the  
 156 cumulative rewards of  $-10 \times 6$  as found in the reward system. The agent had just began favoring  
 157 the more optimal solution of touching its head to reach the goal node before the limit to the number  
 158 of episodes was reached. We hypothesize that the agent would continue to learn the map after  
 159 significantly more episodes to reach the more optimal goal state of the diamond block due to the  
 160 nature of the agent having many more states to explore than those of the previous maps.

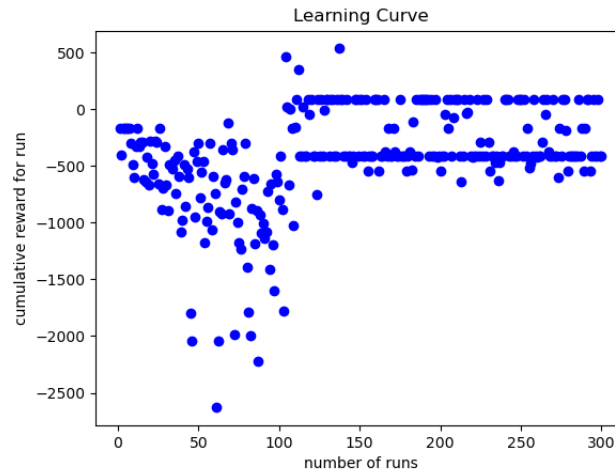


Figure 5: Map 4 Plot

## References

- [1] Choudhary, Ankit, et al. "Introduction to Deep Q-Learning for Reinforcement Learning (in Python)." Analytics Vidhya, 6 May 2019, [www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/](http://www.analyticsvidhya.com/blog/2019/04/introduction-deep-q-learning-python/).
- [2] Ossenkopp, Philip, et al. "Reinforcement Learning – Part 1: Introduction to Q-Learning." Novatec, 13 Feb. 2020, [www.novatec-gmbh.de/en/blog/introduction-to-q-learning/](http://www.novatec-gmbh.de/en/blog/introduction-to-q-learning/).
- [3] Violante, Andre. "Simple Reinforcement Learning: Q-Learning." Medium, Towards Data Science, 1 July 2019, [towardsdatascience.com/simple-reinforcement-learning-q-learning-fcddc4b6fe56](https://towardsdatascience.com/simple-reinforcement-learning-q-learning-fcddc4b6fe56).