
Using Machine Learning to Train an Agent to Terminate Mobs in Minecraft

Farhib S.

Georgia Institute of Technology
Atlanta, GA 30332

Rishav B.

Georgia Institute of Technology
Atlanta, GA 30332

Dominik K.

Georgia Institute of Technology
Atlanta, GA 30332

Jean M.

Georgia Institute of Technology
Atlanta, GA 30332

Saagar M.

Georgia Institute of Technology
Atlanta, GA 30332

Abstract

As robots become more prevalent in society and the workplace, the need for robust algorithms that can learn to control autonomous agents in a wide range of situations becomes paramount. Prior work has shown that deep reinforcement learning models perform reliably in 3D environments, even when rewards are sparse and only visual input is provided to the agent. Here, we use Project Malmö, an AI research platform based on the popular sandbox video game Minecraft, to train an RL agent to combat in-game entities known as “mobs.” We implement a RL algorithm called a Deep Q Network (DQN), as well as a pretrained residual neural network as a baseline model, and compare the differences in performance. We expect that with minimal hyperparameter tuning, the RL model will learn significantly more than the baseline, and that the agent will succeed in defending itself to some extent.

Introduction and Background

Minecraft is a popular sandbox video game that contains a number of hostile non-player entities known as “mobs”; these entities are meant to attack and kill the player character. Our agent will have to learn strategies to deal with each type of hostile mob with the goal of defeating as many mobs and surviving as long as possible. Additionally, the environment in a Minecraft “world” can be randomly generated using an algorithm or built by the player. To create a closed environment for our agent to learn and fight against these mobs, we will be using Microsoft’s Project Malmo. Using machine learning in minecraft is the focus of a large competition called MineRL, which provides rigorous guidelines towards achieving an agent that can operate autonomously in minecraft. It is our hope that methods like the ones we are using to train an agent in a simulated environment can be extrapolated to real life applications like robotics in the physical world. Since minecraft as an environment is completely customizable, it makes it ideal for entry level testing of potential real world use cases.

Problem Definition

The agent will have to last as long as possible while defeating as many distinct hostile entities as possible and navigating the environment. The agent will receive positive rewards for defeating

28 entities/surviving and negative rewards for being defeated and losing health itself. We are utilizing a
29 fairly dense reward structure, with the hope that this will enable the agent to learn good behaviors
30 more reliably. Since we are rewarding the agent for successful hits on mobs and survival, and are
31 negatively rewarding it for taking damage and dying, we can see our reward structure is dense.
32 Additionally, to increase the chance of the agent learning the reward for attacking mobs, we let the
33 agent continually attack, so it has to learn to face the mobs, rather than face them and then attack.
34 Below are listed the present actions and rewards we used to train our preliminary RL model:

- 35 • Action Space: Move Forward, Move Backward, Turn Left, Turn Right, Do Nothing
36 • Rewards: Rewards: Death (-100), Damage Taken (-4), Damaging Zombie (15), Per Action
37 Taken (0.05), Zombie Killed (150)

38 Data Collection

39 Since we are using Deep Q Learning, we did not have to collect any data. The agent's observations in
40 the environment was our "data," on which the neural network trained on. The agents observations in
41 the environment were 640x480, which we rescaled to an 84x84 image.

42 Definitions

43 **Step:** A step is every iteration in an episode. Each step, the agent makes an observation, takes an
44 action, and learns from previous memories.
45 **Episode:** Each run of the game in which the agent plays (until it dies) is called an episode. Reward:
46 The agent receives a positive reward for being in a good state and taking an optimal action like hitting
47 a zombie. It receives a negative reward for things like getting hit.
48 **Q-value:** The Q value is essentially a numeric value assigned to a state action pair determining how
49 "good" that action is given the current state.
50 **DQN:** A neural network that we are using to approximate the Q-value of a state action pair. Target
51 Network: It is a copy of the DQN, but is only updated periodically. This is used to increase the
52 stability of the algorithm.
53 **ResNet50:** A large image recognition CNN.

54 Methods

55 We used a Convolutional Deep Q Network to take in the image input and output what action(s) to
56 take. One of the ways that Project Malmo allowed our agent to "see" in the Minecraft world was
57 through images, so using a convolutional neural network made logical sense. Similar to most CNNs,
58 we started with the CNN workflow (Convolution, Max Pooling, Activation) and then used some
59 fully connected layers. We also used a replay buffer to allow the agent to have "memory," giving
60 the agent a way to utilize past trials. Another implementation detail is that we used a target network
61 that we copied the weights to periodically every (300 steps) so that our DQN converged to a more
62 stable solution. As prior research had shown us, using a recurrent neural network would not give us
63 significant improvements so this is not a path we decided to follow [2].
64 With those implementation details, we followed the regular Q-learning algorithm, which is as follows:

- 65 1. Get state of the environment
- 66 2. Take action using an epsilon greedy policy
- 67 3. Create SARS tuple (state, action, reward, best state) and store it into replay buffer
- 68 4. Sample from the replay buffer and use the sample to update the weights of the DQN.
- 69 5. Update target network (only do this once every 300 episodes).

70 As a baseline model we took the feature representation from a large pre-trained CNN such as
71 ResNet50, by using the model and excluding the final dense layer, and using this in place of our
72 convolution layers. We had predicted that this would likely get us some performance, but would
73 inherently be worse, since we had fixed some of our trainable parameters.
74 Using a DQN on Minecraft is not a very novel method as Minecraft is a fully-observable, deterministic
75 environment, which is well suited for reinforcement learning.

76 **Metrics**

77 While we do have a loss function we can plot to track learning, for our reinforcement learning
78 problem, tracking the metric of total reward per episode is a better measure of our progress. This
79 metric is essentially how well the agent played in its environment during that specific episode. We
80 do not really have data to split for K-fold validation, but we have trained entirely new models for
81 different configurations of our hyperparameters and evaluate them based on the total reward per
82 episode metric. Each of these models are trained on new “data” since each run of each episode will
83 be different. The loss function we used for the DQN is MSE (mean-squared error) since our Q-value
84 function is a continuous function and approximating the Q-value function for a continuous input state
85 is a regression problem. Applying gradient descent or another optimization method to minimize this
86 loss function allowed our network to learn the Q-value function.

87 **Initial Results**

88 The results of our training did show learning within our reward scheme, but that reward scheme was
89 not optimal for what we wanted our agent to learn. With the large negative reward of -1000 for dying
90 and the maximum number of steps set to 50, the agents reward was usually either 50 (since it survived
91 the entire episode) or something less than -950. Because of this, any other rewards that could have
92 been explored would be overlooked since the negative reward for dying had such a large magnitude.
93 We can see this in the graph below, in which the rewards are sporadic and shifting between around 0
94 and -1000. Although as time passes, we see that the density of rewards that were around the 0 mark
95 increases while the opposite occurred for the -1000 rewards. The Savitzky-Golay smoothing filter
96 visualizes this nicely. While these results were less than ideal, we were still able to get something
running and learning in the minecraft environment, which was our main goal for this touchpoint.

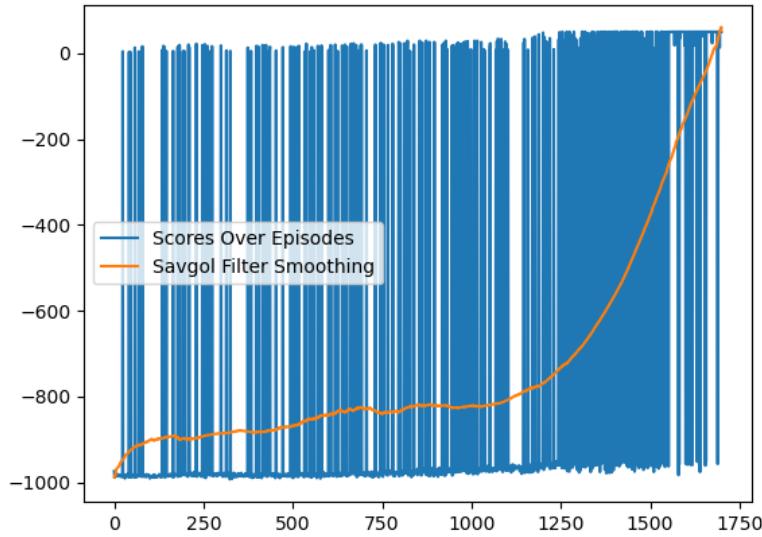


Figure 1: Initial agent rewards while training our DQN

97
98 While neither of the networks performed optimally, the CNN did perform better than ResNet50 as
99 expected.

100 **Final Results**

101 As we had expected, the baseline performed worse for a couple of reasons. Because the feature
102 representation of ResNet50 is so large, this model ended up having more parameters to train even

103 though the dense layers that we added were of the same size. This made it a lot more challenging to
 104 train: the main reason the baseline did not perform as well as the CNN was because the overhead
 105 required for a backward pass through the larger dense layers caused the agent to take actions at a
 106 much slower rate. Because of this, instead of taking actions a few times each second it took a few
 107 seconds per action. This knocked performance down by a lot since the chances of the agent actually
 108 hitting the zombies decreased by a large factor. This makes killing zombies a lot more challenging
 109 since the agent must line its sword as the agent takes the attack action (which now happens less
 110 frequently). Although there are many hindering factors, the baseline agent still was able to show
 111 some initial learning, but the overall rewards were significantly lower than the regular agent and the
 112 algorithm did not converge. Below is the graph for baseline training.

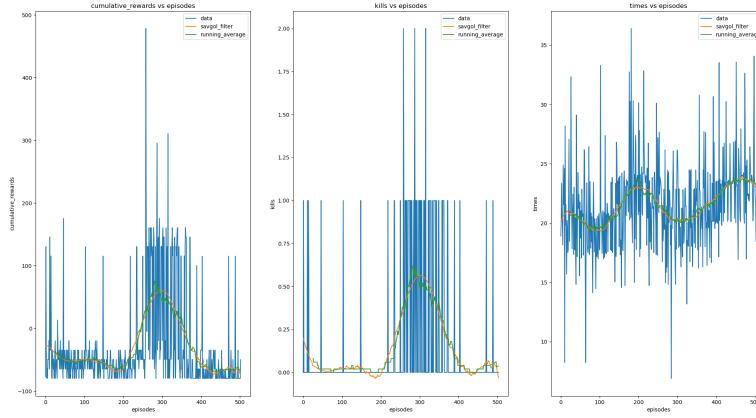


Figure 2: Reward, Kills, and Time Survived per Episode for ResNet50

113
 114 As for the CNN model, we were able to show learning by the agent as shown in the graphs below.

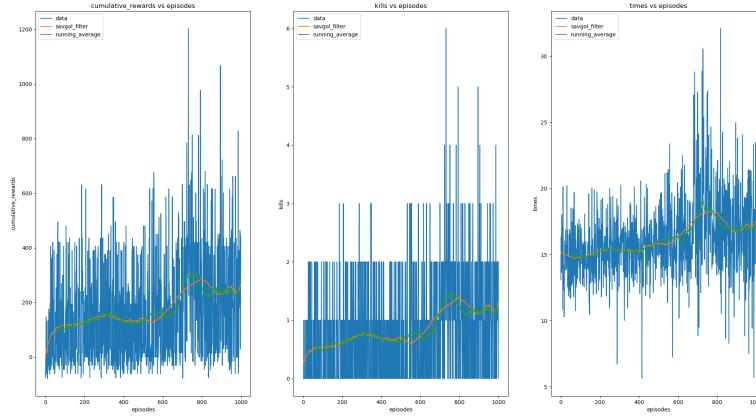


Figure 3: Reward, Kills, and Time Survived per Episode for DQN

115
 116 However, the strategies of the agent were not optimal. The main strategy the agent followed was to
 117 put itself into a corner and then spin or just face the zombies. While this does allow the agent to
 118 hit the zombies, it also allows the zombies to hit the agent, which is suboptimal. One of the better
 119 strategies the agent learned was that if the zombies' locations were unknown but the agent was getting

120 hit, it would face a wall and then back up, which usually put the zombies back in view as they would
 121 follow. This also allowed the agent to get a few hits on the zombie. The baseline agent ended up
 122 learning similar strategies but implemented them worse because of the hindering factors previously
 123 mentioned.



Figure 4: A Screen Capture of the Agent Learning to Back Away

124
 125 Overall, the performance of the agent was okay, but it had room for improvement. Unfortunately, we
 126 did not get enough time to tune the hyperparameters and reward scheme as much as we had liked,
 127 due to limitations in time and compute. Due to the fact that we were training this on a laptop, each
 128 train of the agent took around 6-10 hours. So, we were only able to train a few different versions
 129 of the model, our best of which ended being version 3. Below is one of the earlier versions of the
 130 model, which did not learn very well since the agent was too scared to try anything since the penalty
 131 for getting hit was too high.

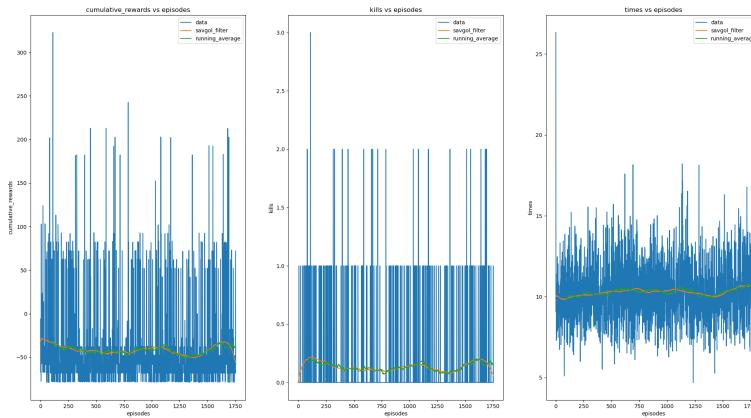


Figure 5: Reward, Kills, and Time Survived per Episode for DQN with Poor Hyperparameters

132

133 **Conclusion**

134 In this project we learned more about Deep Q networks and practical convolutional neural networks
135 work and utilized them. To improve upon our research we would add things like current health,
136 zombie location, and ambient noises to the observation space of the agent in order to alleviate the
137 problem of the agent taking nearly random actions when zombies are not on the screen. These are
138 values that we can usually deduce while playing but the agent had no way of directly seeing these
139 values. Given more time we could further tune hyperparameters, and rewards, hopefully letting the
140 agent learn optimal behavior more reliably and quickly. Also we would have used an auto encoder
141 instead of ResNet50, as we could have avoided many of the problems with processing time taking
142 too long with that approach.

143 **Ethics Statement**

144 The societal impacts of this project are difficult to determine due to the fact that the agent
145 works with a videogame. The most direct impact of this project would be to improve video
146 game artificial intelligence, thus improving the experience of playing video games. However,
147 if we were to consider our project as an agent performing actions in response to visual stimuli,
148 we can expand the potential societal impacts to a much wider set of applications, especially in robotics.

149
150 Because of the wide range of applications robotics brings, the societal benefits and harms
151 are equally numerous. Ensuring that the benefits sufficiently outweigh the harms comes down to two
152 main things: the accuracy of the model itself and the quality of the data. The model presented here
153 had some accuracy issues that could result in significant problems if applied to crucial services such
154 as surgery or medical care in general. The solution to this issue would be to tune hyperparameters
155 until the model performed acceptably.

156
157 Quality of data is arguably even more important than the model itself; it determines the bi-
158 ases and features the model will focus on. Special thought needs to be taken to ensure quality of
159 input data the model will train on to ensure undue. If our agent was applied to trying to identify
160 human faces, for example, this would take the form of diversity in the dataset to allow for better
161 identification of different races. Ignoring these consideration could potentially result in significant
162 societal harm.

163 **References**

- 164 [1] Christian S., Yanick S. & Manfred V. (2020). Sample Efficient Reinforcement Learning through Learning
165 from Demonstrations in Minecraft. arXiv. Retrieved March 1, 2021, from <https://arxiv.org/abs/2003.06066>
166 [2] Clément R. & Vincent B. (2019). Deep Recurrent Q-Learning vs Deep Q-Learning on a simple
167 Partially Observable Markov Decision Process with Minecraft. arXiv. Retrieved March 1, 2021, from
168 <https://arxiv.org/abs/1903.04311>
169 [3] Volodymyr M., Koray, K., David, S., Alex, G., Ioannis A., Daan W. & Martin R. (2013). Playing Atari with
170 Deep Reinforcement Learning. arXiv. Retrieved March 1, 2021, from <https://arxiv.org/abs/1312.5602>
171
172 Github Pages Link: <https://minerl.bhagat.io/>