

# Drone Flight Controller, Project 9 - MLiS Part I

Jarrad Foley<sup>1,\*</sup> and Stephen McSweeney<sup>1,†</sup>

<sup>1</sup>*School of Physics and Astronomy, University of Nottingham, Nottingham, NG7 2RD, UK*

(Dated: January 19, 2023)

This paper outlines the process and intricacies that were explored with creating a drone that is controlled autonomously through reinforcement methods. The method implemented is a state-action value algorithm with a soft policy[1] to ensure all states have a non-zero probability of being taken. This is kept to the bounds of a defined space and set of targets and a specified amount of timesteps. The data demonstrated a logarithmic improvement in the drones performance over epochs. In our sample of completed simulations, it performed better than the heuristic model after an average of 2596 Epochs. Further simulation is required for statistical significance in demonstrating better performance than the Heuristic model.

## I. INTRODUCTION

Drones have become a useful tool in the modern era, with many applications such as agriculture, photography and warfare. Defining movement of a drone is a non-trivial task with many variables to consider. While there are many ways to control a drone such as user controlled, this project focuses on implementing autonomous flight for a 2d simulated drone. This drone will fly within a defined space to hit a number of predefined targets. The drone is controlled using two motors which define the thrust independently allowing for navigation, and is subject to gravity and drag forces on rotation and planar motion. The task is to write a flight controller which uses reinforcement learning to have the drone traverse space to hit as many targets as it can within a specified amount of time.

## II. IMPLEMENTATION

For the implementation for this project we used a Soft policy Monti Carlo algorithm implementing a state-action value function ensuring that all actions can be taken with a non-zero probability.

### A. States

The states for this model uses the pitch of the drone, distance from the target both x and y, and velocity. The State space used is  $s = \{\Theta, \Delta X, \Delta Y, V\}$  Where each variable is discretized such that: pitch can take 30 values, X can take 40 values, Y can take 30 values, and Velocity can take 10 values. The total number of possible states is 360000.

### B. Actions

There are 6 actions used within this program to define what movement the drone does.

$$Ta = \frac{1}{10} \max(-1, \min(10\Delta y, 1))$$

The first action is used to counteract the drones current pitch if it is facing the wrong direction from the target.

$$T_1, T_2 = \begin{cases} T_1 = \frac{1}{2}, T_2 = \frac{1}{2} \\ \text{where } \text{sgn}(\Delta x) = \text{sgn}(\Theta) \\ T_1 = \max(1, \min(\frac{\text{sgn}(\Delta x)}{2} + 1), 0), \\ T_2 = \max(1, \min(-(\frac{\text{sgn}(\Delta x)}{2} + 1), 0)) \\ \text{where } \text{sgn}(\Delta x) \neq \text{sgn}(\Theta) \end{cases}$$

The second action adjusts the thrusts based on where the target is in the space and the angle of change that the drone will need to adjust to be directed at the target multiplied by 2. If the target is to the left of the drone it will have a lower left thrust than right thrust and vice versa.

$$T_1 = \max(1, \min(\frac{2}{5} + Ta + 2\Delta\text{pitch}, 0)) \\ T_2 = \max(1, \min(\frac{2}{5} + Ta - 2\Delta\text{pitch}, 0))$$

The third action gets the thrust based on the velocity of the drone's pitch to adjust which thrust it needs more power to counteract the current rotation when it's not rotating towards the target. This action is unstable compared to another action which does the same thing just at a slower rate.

$$T_1 = \frac{1}{2} - \frac{2}{5} \text{sgn}(\text{pitchvelocity}) \\ T_2 = \frac{1}{2} + \frac{2}{5} \text{sgn}(\text{pitchvelocity})$$

The fourth action changes the thrust depending on the current pitch velocity of the drone. If the pitch velocity is in the direction of the target then the thrusts will more equal to slow down how much it spins, however if the pitch velocity is going in the opposite direction of the target then it strongly sets the thrusts to alter the pitch velocity to be towards the target, spinning it into the correct direction.

---

\* [ppxf2@nottingham.ac.uk](mailto:ppxf2@nottingham.ac.uk)

† [ppxsm5@nottingham.ac.uk](mailto:ppxsm5@nottingham.ac.uk)

$$T_1, T_2 = \begin{cases} T_1 = \max(\frac{2}{3}, \min(\text{sgn}(\Delta x), \frac{3}{5})) \\ T_2 = \max(\frac{3}{5}, \min(\text{sgn}(\Delta x), \frac{2}{3})) \\ \text{where } \text{sgn}(\Delta x) = \text{sgn}(\text{pitchvelocity}) \\ \\ T_1 = \max(\frac{1}{5}, \min(\text{sgn}(\Delta x), \frac{4}{5})) \\ T_2 = \max(\frac{4}{5}, \min(\text{sgn}(\Delta x), \frac{1}{5})) \\ \text{where } \text{sgn}(\Delta x) \neq \text{sgn}(\text{pitchvelocity}) \end{cases}$$

The fifth action adjusts the thrusts based on where the target is in the space and the angle of change that the drone will need to adjust to be directed at the target. If the target is to the left of the drone it will have a lower left thrust than right thrust and vice versa.

$$T_1 = \max(1, \min(\frac{1}{2} + Ta + \Delta\text{pitch}, 0))$$

$$T_2 = \max(1, \min(\frac{1}{2} + Ta - \Delta\text{pitch}, 0))$$

The final action is used to adjust the drones y value. If the drones y value is lower than the target's y value then the drone will ascend and if the drones y value is above the target's y value then the drone will descend.

$$T_1 = \max(0, \min(\text{sgn}(\Delta y), 1))$$

$$T_2 = \max(0, \min(\text{sgn}(\Delta y), 1))$$

### C. Rewards

We use the following reward function:

$$R = \frac{(-3a + 0.1b - c + 2d^3)}{100}$$

Where:

$$a = \|TargetDistance\|$$

$$b = \text{abs}(DroneVelocity)$$

$$c = \lfloor \text{abs}(DronePitch) \rfloor$$

$$d = 1 + TargetsHit$$

## III. RESULTS

Each assessment was done with the following target coordinates.

Target	1	2	3	4	5
x	0.35	-0.35	0.5	-0.35	0.35
y	0.3	0.4	-0.4	0	0.4
Target	6	7	8	9	10
x	-0.15	-0.35	0.35	-0.5	0.35
y	-0.1	-0.3	-0.4	0.4	0

Each simulation was also run for 3000 simulation steps and 5000 epochs.

In this simulation, the heuristic algorithm hit 5 targets and achieved a value for the sum of rewards as 4894.42, against our chosen reward function.

The following graph shows an example of the Soft state Monti-Carlo state machine algorithm's results, under the same conditions. Figure 1.

We modelled the cumulative reward function for any given epoch using the least-squares optimized log function. Figure 1. The implementation is not deterministic

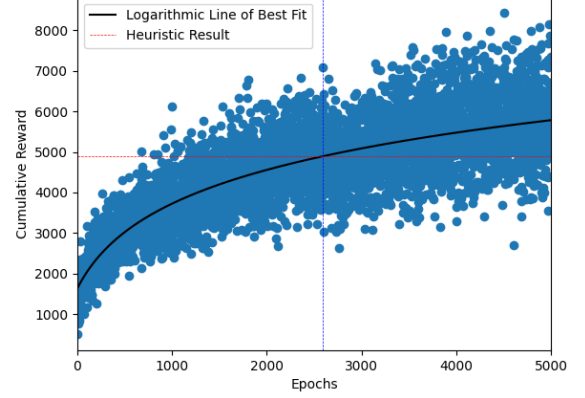


FIG. 1. Shows Logarithmic line of best fit of the data with result of the heuristic.

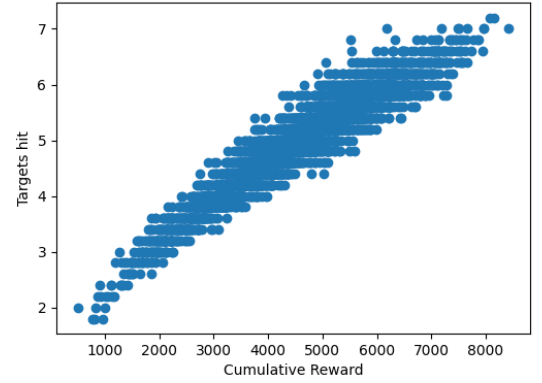


FIG. 2. Shows correlation between how many targets hit and cumulative reward for simulation.

therefore an average of five simulations is used to reduce the error in our least squares exponential model.

The Logarithmic line for best fit is defined by the following equation:

$$1475.01044 * \ln(0.00310817391CumulativeReward) + 1642.65857$$

The model has a residual standard deviation of 711.29847146.

Figure 2 shows targets hit against cumulative rewards.

## IV. EVALUATION

This graph demonstrates strong correlation between the cumulative reward function and the number of targets hit suggesting it is a good estimator for the performance of the drone. Figure 1

The best fit logarithmic model suggests that we would expect the simulation to overtake the heuristic model with its average simulation after 2596 epochs. Figure 2

However, the model has a large Mean Squared Error due to only averaging over 5 simulations. The model has a residual standard deviation of 711.29847146, and we would therefore require the cumulative reward to be 7028 for the heuristic model's result to lie outside of 3 standard deviations. The model therefore suggests that after 12395 epochs we would expect any results worse than the heuristic algorithm to be anomalies. Due to our model only being trained to 5000 epochs, however, the model may struggle to be representative of the data at this point.

## V. CONCLUSIONS

The ML algorithm successfully demonstrated logarithmic improvement over time, and therefore our actions

and state space were appropriate for achieving some level of training. Also, the cumulative reward function was a successful indicator for the performance of the drone to the high correlation. The model shows signs of outperforming the Heuristic model but more simulations are required to demonstrate this with statistical significance. This was a limit of the time constraint of the project.

For future improvements to this work, more simulations are recommended and over more epochs to improve the model. Through collecting data for different actions and state spaces, it may be possible to improve the performance of the model.

- 
- [1] R. S. Sutton and A. G. Barto, [Reinforcement learning: An introduction](#). (2018).