

Bachelor Thesis

# Reinforcement Learning for Control of Flying Drones

Johannes M. Tölle

Date of Submission: 23. 08. 2023

Advisor: Prof. Dr. Carlo D'Eramo



Julius-Maximilians-Universität Würzburg  
Lehrstuhl für Informatik VII  
Robotik und Telematik

# Declaration

I certify that I have authored this thesis, including all attached materials, constantly and without the use of any other aids than those specified.

All passages taken literally or analogously from published or unpublished works are clearly identified as such in each individual case, stating the source.

The work has not yet been submitted in the same or a similar form as a thesis.

I am aware that violations of this declaration and deliberate deception may lead to a grade of 5.0.

Würzburg, June 17, 2023

---

Johannes M. Tölle

# Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>2</b>  |
| <b>2</b> | <b>Basics</b>  | <b>3</b>  |
| 2.1      | Quadrocopter Basics . . . . .                                  | 4         |
| 2.1.1    | Quadrocopter Flight Dynamics . . . . .                         | 5         |
| 2.1.2    | Autonomous Quadrocopter . . . . .                              | 8         |
| 2.2      | Reinforcement Learning Basics . . . . .                        | 11        |
| 2.2.1    | Classification of Machine Learning Areas . . . . .             | 12        |
| 2.2.2    | Markov Decision Process . . . . .                              | 14        |
| 2.2.3    | V,Q Function & The Bellman Equations . . . . .                 | 16        |
| 2.2.4    | Neural Networks . . . . .                                      | 20        |
| 2.2.5    | State of the Art Algorithms . . . . .                          | 23        |
| 2.3      | Pybullet Physics Simulator & Gym Pybullet Drones . . . . .     | 26        |
| 2.3.1    | BaseAviary Class . . . . .                                     | 26        |
| 2.3.2    | BaseSingleAgentAviary Class . . . . .                          | 26        |
| <b>3</b> | <b>Flight-Control Concept</b>                                  | <b>27</b> |
| <b>4</b> | <b>Implementation</b>  | <b>28</b> |
| 4.1      | Environment Classes . . . . .                                  | 29        |
| 4.1.1    | WindSingleAgentAviary Environment Class . . . . .              | 29        |
| 4.1.2    | WindSingleAgentPathfollowingAviary Environment Class . . . . . | 36        |
| 4.2      | Scripts & Evaluation Tools . . . . .                           | 37        |
| 4.2.1    | Learning Script . . . . .                                      | 37        |
| 4.2.2    | Evaluation Script . . . . .                                    | 38        |
| 4.2.3    | Evaluation Tools . . . . .                                     | 39        |
| <b>5</b> | <b>Evaluation</b>  | <b>41</b> |
| 5.1      | Drone Model . . . . .  | 41        |
| 5.2      | Setup . . . . .  | 42        |
| 5.3      | Metric . . . . .   | 43        |
| 5.4      | Results . . . . .  | 44        |
| <b>6</b> | <b>Conclusion &amp; Future Work</b>                            | <b>45</b> |
| <b>7</b> | <b>Appendix</b>  | <b>46</b> |
|          | <b>Bibliography</b>  | <b>47</b> |

# List of Figures

|      |   |    |
|------|---|----|
| 2.1  | "X" and "+" configuration of a quadcopter with the motors ( $M_1...M_4$ ), the flight controller( $FC$ ) and the rotational movement of the corresponding propeller (red clockwise, blue counterclockwise) . . . . .  | 4  |
| 2.2  | Visualization of Roll, Pitch and Yaw rotational movement (from left to right) and the corresponding difference in the rotational speed $\omega_i$ . . . . .   | 6  |
| 2.3  | Sketch of an autonomous quadcopter flight control system . . . . .  | 8  |
| 2.4  | Concept onion of Artificial Intelligence [DJS20] . . . . .  | 11 |
| 2.5  | The main learning strategies of machine learning including the related approaches like classification and regression under the Supervised Learning, clustering under the Unsupervised Learning and a sketch of the Reinforcement Learning concept . . . . .   | 12 |
| 2.6  | The basic concept of RL: the agent takes an action based on the observed state and receives a numerical reward at each time step . . . . .  | 13 |
| 2.7  | Visualization of a simple MDP with 4 states ( $s_0, s_1, s_2, s_3$ ) (blue nodes), 2 actions ( $a_0, a_1$ ), state-action edges (green) and action-state edges (dark red) containing a tuple (possibility, reward) . . . . .  | 15 |
| 2.8  | The McCulloch-Pitts model of a single neuron uses a weighted sum with the inputs $x_i$ and weights $w_i$ and a non-linear activation function $g()$ in order to compute the output $z$ [Bis94] . . . . .  | 20 |
| 2.9  | 8 different Activation Functions. . . . .   | 21 |
| 2.10 | Simple Feed-Forward Network with 3 neurons in the input layer, 2 hidden layers with 4 neurons and with 2 layers in the output layer . . . . .   | 22 |
| 2.11 | Overview of the different categories of RL algorithms with examples: PPO and DDPG Algorithm coloured because they were used in this work. . . .   | 23 |
| 3.1  | Proposed Concept of autonomous flight control with the use of a intelligent agent implemented with the use of a NN. The output of the NN is denormalized in order to translate them to motor signals. With the use of the sensors and position estimation the input for the next control step are calculated. . . . . | 27 |
| 4.1  | Concept of the implemented software: the different tools(red), scripts(violet) that are used in order to learn the intelligent Agent robust flight control with the use of RL. . . . .  | 28 |
| 4.2  | Visualization of the used reward function with the goal (0,0,0.5) and a colour scale: the distance dependent part (a) and the big penatalization, because the drone hit the ground (b) . . . . .  | 34 |

|     |   |    |
|-----|---|----|
| 7.1 | shortened excerpt of the python implementation of the learning script . . . | 46 |
|-----|---|----|

## List of Tables

|     |   |    |
|-----|---|----|
| 2.1 | Example of all two step trajectories with the matching probability, reward and expected reward of the MDP shown in Figure 2.7 . . . . . | 17 |
| 4.1 | The different ActionTypes with the corresponding dimensionality of the action, its range and how it is processed. . . . .               | 33 |

# List of Algorithms

|   |  |    |
|---|--|----|
| 1 | Proximal Policy Optimization [SWD <sup>+</sup> 17] | 25 |
| 2 | Learning Script                                    | 37 |
| 3 | Evaluation Script                                  | 38 |

# Abstract

content...



# 1 Introduction

## 2 Basics

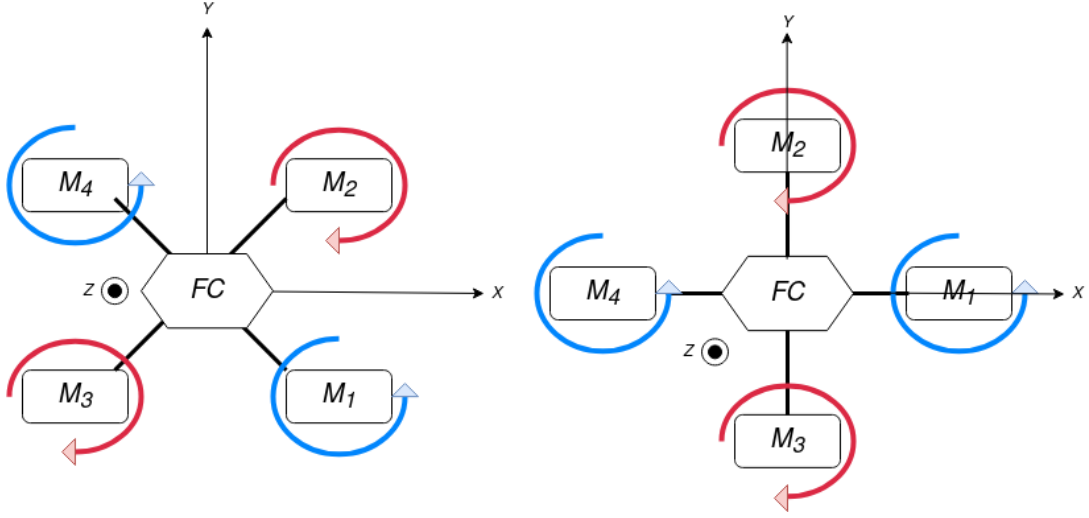
The task to perform autonomous Drone-Control may be described as a rather challenging task, because it combines a couple of different topics from various areas of computer science:

*Quadrocopter* possess rather unique flight dynamics and their flight control is its own unique research topic with large requirements. In order to achieve accurate flight control it is crucial to use sensor data with a high accuracy, process and apply them in realtime to the actuators. At the same time there is research to improve the error tolerance of flight control and make it more adaptive to a change in the environment or the payload. These Errors can occur due to turbulent conditions like wind.

*Reinforcement Learning* is a natural approach that is known to show more accurate results than the classic PID Controller in many use cases, because of the ability to approximate and generalize decision making. In Addition, a control implemented with an intelligent agent that had been learned with RL is often more robust to environment changes and is capable of adaption to these changes.

Although training a real drone from the beginning on flight seems more intuitive, this shows some risks. The first learning episodes show a bad behaviour that may risk the real hardware of the drone. A Reinforcement Learning episode is in need of quite complex calculations that would overstrain the hardware. Also, there can be multiple parallel training episodes and simulation time must not be synced to the real world time. As a consequence of this the time the learning takes can be decreased at the cost of complexity to translate the simulation to the real world.

This chapter provides the needed basics in order to understand the idea of using RL for robust Drone-Control. Therefore, it provides a summary of the quadrocopter flight dynamics and sketches the autonomous quadrocopter concept. In addition essential parts of conventional autonomous quadrocopter are explained. Then, the basic concept of Markov Decision Processes and RL are presented. At last, the used Physics Simulator is explained as well as the extension that can be used for RL quadrocopter problems.



**Fig. 2.1:** "X" and "+" configuration of a quadcopter with the motors ( $M_1...M_4$ ), the flight controller ( $FC$ ) and the rotational movement of the corresponding propeller (red clockwise, blue counterclockwise)

## 2.1 Quadcopter Basics

A quadcopter is a modern aircraft with a wide range of applications which is immensely popular in modern society. Typically, a quadcopter can be described as an aircraft with 6 degrees of freedom (DOF), which consist of 3 translational and 3 rotational DOF around the  $x, y, z$  axis with the corresponding angle  $\Theta, \phi, \psi$ .

As a consequence it has the ability to maneuver in 3D space with all possible rotations, although some rotations like a rotation of  $\pi$  around the  $x$  Axis ( $\Theta = \pi$ ) may cause an inevitable crash.

The heart of the quadcopter is the flight controller ( $FC$ ) which consists of sensors and a MCU that controls the flight. Besides, the aircraft possesses 4 motors with corresponding propellers. There are different configurations how the motors can be placed in relation to the axis. The two most common configurations, the "X" and "+" configuration, can be seen in the above figure. The "+" configuration places the motors alongside the  $x$  and  $y$  axis. In contrast to that each motor has an angle of  $45^\circ$  to the  $x$  and  $y$  axis in the "X" configuration.

Each of the propellers can be described by their rotational speed  $\omega_i$  with  $i \in [1...M]$  and spins either clockwise or counterclockwise (Figure 2.1). Depending on the configuration of the motors, the rotation speeds directly influence the *Euler Angles*  $\Theta, \phi, \psi$  and the upward thrust  $f$ . However, hereafter the most common "X" configuration is used to explain the quadcopter further.

### 2.1.1 Quadcopter Flight Dynamics

The quadcopter flight dynamics mainly depend on the rotational speed  $\omega_i$  of the propellers, which mainly produce the translational and rotational movements. Besides, the thrust factor  $b$  which is a constant that mainly consists of propeller geometry and frame characteristics also influences the movement. Depending on the different  $\omega_i$ ,  $b$  an upward thrust  $f$  and a rotation  $u_\Theta, u_\phi, u_\psi$  is produced that can be comprehended with the following set of formulas:

$$u_f = b \cdot \sum_{i=1}^4 \omega_i^2 = b \cdot (\omega_1^2 + \omega_2^2 + \omega_3^2 + \omega_4^2) \quad (2.1)$$

$$u_\Theta = b \cdot (\omega_1^2 - \omega_2^2 + \omega_3^2 - \omega_4^2) \quad (2.2)$$

$$u_\phi = b \cdot (\omega_1^2 + \omega_2^2 - \omega_3^2 - \omega_4^2) \quad (2.3)$$

$$u_\psi = b \cdot (\omega_1^2 - \omega_2^2 - \omega_3^2 + \omega_4^2) \quad (2.4)$$

#### Hover, Rise & Fall

*Hovering* can be described as holding a pose  $p = x_p, y_p, z_p$  above the ground, rising can be seen as increasing  $z_p$  and falling as decreasing  $z_p$ . Therefore, the quadcopter should not have a rotation on  $x, y, z$  axis and mainly depends on Equation (2.1).

The product of the thrust factor  $b$  and the sum of the squared rotational speeds produces an upward thrust  $f$ . In order to hover equal rotational speeds must be applied to each motor/propeller in order to avoid a rotational movement, because:

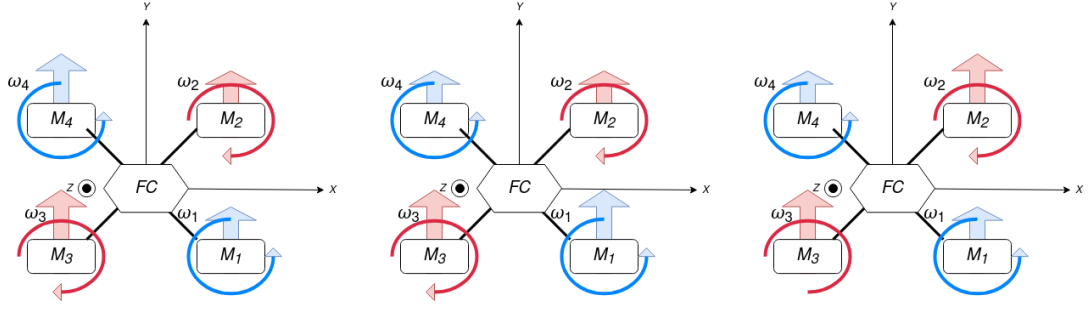
$$u_\Theta = b \cdot (\omega_1^2 - \omega_2^2 + \omega_3^2 - \omega_4^2) = b \cdot 0 = 0 \quad (2.5)$$

$$u_\phi = b \cdot (\omega_1^2 + \omega_2^2 - \omega_3^2 - \omega_4^2) = b \cdot 0 = 0 \quad (2.6)$$

$$u_\psi = b \cdot (\omega_1^2 - \omega_2^2 - \omega_3^2 + \omega_4^2) = b \cdot 0 = 0 \quad (2.7)$$

Since all rotational movements null itself, there is only an upward thrust as aerodynamic effect. Depending on the upward force  $F$  and gravitational Force  $g$  the thrust produces the already mentioned movements:

- $F < g$ : Fall
- $F = g$ : Hover
- $F > g$ : Rising



**Fig. 2.2:** Visualization of Roll, Pitch and Yaw rotational movement (from left to right) and the corresponding difference in the rotational speed  $\omega_i$

## Roll

In order to *roll* (rotation around  $y$  axis) the thrust on one side has to be greater than on the other side. For example, increasing  $\omega_3, \omega_4$  would lead to uneven distributed thrust. The left side of the aircraft experiences more thrust (Figure 2.2) and as a consequence the aircraft experiences a torque about the  $y$  axis. Equation (2.2) to Equation (2.4) also shows this aerodynamic effect:

$$\begin{aligned}
 \omega_3^* &= \omega_4^* = n \cdot \omega_1 = n \cdot \omega_2 & n > 1 \\
 u_\Theta &= b \cdot ((\omega_1)^2 - (\omega_2)^2 + (\omega_3^*)^2 - (\omega_4^*)^2) = 0 \\
 u_\phi &= b \cdot ((\omega_1)^2 + (\omega_2)^2 - (\omega_3^*)^2 - (\omega_4^*)^2) < 0 \\
 u_\psi &= b \cdot ((\omega_1)^2 - (\omega_2)^2 - (\omega_3^*)^2 + (\omega_4^*)^2) = 0 \\
 u_f &= b \cdot ((\omega_1)^2 + (\omega_2)^2 + (\omega_3^*)^2 + (\omega_4^*)^2) > 0
 \end{aligned}$$

Analogue a roll to the other side can be accomplished by increasing  $\omega_1, \omega_2$ , then  $u_\phi > 0$ . Since the pairs  $(\omega_1, \omega_2)$  and  $(\omega_3, \omega_4)$  each have one propeller that spins clockwise and one propeller that spins counterclockwise, the total torque of the quadrocopter does not change.

Equation (2.1) shows that there still is an upward thrust, but due to the the rotation of the aircraft the thrust force  $F$  has a component in  $x$  direction. This results in a translational movement along the  $x$  axis.

## Pitch

A *pitch* movement around  $x$  axis works similar to a roll movement, just with other motor pairs. In this case, the thrust at the front or at the back has to increase. For example, increasing  $\omega_3, \omega_1$  would lead to a bigger thrust at the back and the corresponding torque (Figure 2.2). Equation (2.2) to Equation (2.4) also shows this aerodynamic effect:

$$\begin{aligned}\omega_3^* &= \omega_1^* = n \cdot \omega_4 = n \cdot \omega_2 & n > 1 \\ u_\Theta &= b \cdot ((\omega_1^*)^2 - (\omega_2)^2 + (\omega_3^*)^2 - (\omega_4)^2) > 0 \\ u_\phi &= b \cdot ((\omega_1^*)^2 + (\omega_2)^2 - (\omega_3^*)^2 - (\omega_4)^2) = 0 \\ u_\psi &= b \cdot ((\omega_1^*)^2 - (\omega_2)^2 - (\omega_3^*)^2 + (\omega_4)^2) = 0 \\ u_f &= b \cdot ((\omega_1^*)^2 + (\omega_2)^2 + (\omega_3^*)^2 + (\omega_4)^2) > 0\end{aligned}$$

Analogue a roll to the other side can be done by increasing  $\omega_2, \omega_4$ , then  $u_\Theta < 0$ . The total torque also stays the same and the thrust force  $F$  gets a component in  $y$  direction. This results in a translational movement along the  $y$  axis.

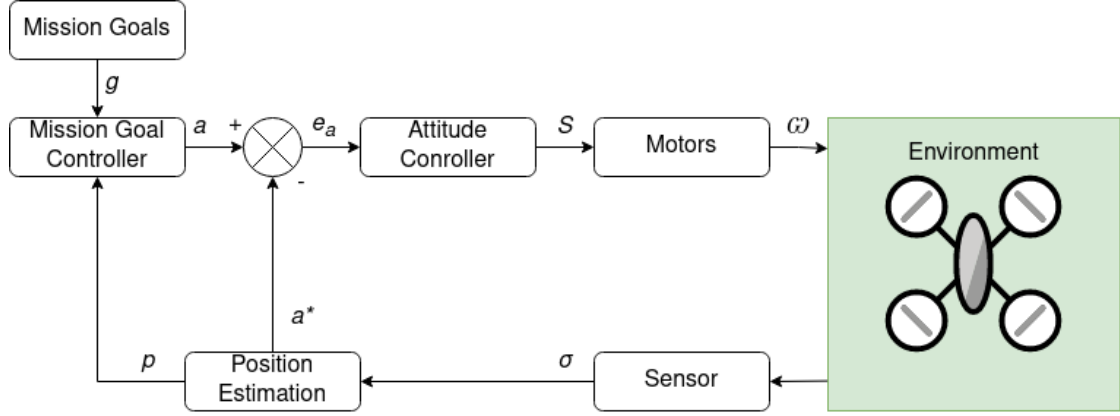
## Yaw

A *yaw* movement around  $z$  axis is not implemented by using a difference in thrust like in the other rotational movements. Furthermore, differences in torque are used to accomplish yaw rotations (Figure 2.2).

For example, if  $\omega_1, \omega_4$  increase, the quadcopter rotates clockwise, because the sum of all torques has to stay the same. The external yaw rotation compensates the difference of the sum of the torques of the propeller. Besides, Equation (2.2) to Equation (2.4) shows this aerodynamic effect:

$$\begin{aligned}\omega_4^* &= \omega_1^* = n \cdot \omega_3 = n \cdot \omega_2 & n > 1 \\ u_\Theta &= b \cdot ((\omega_1^*)^2 - (\omega_2)^2 + (\omega_3)^2 - (\omega_4^*)^2) = 0 \\ u_\phi &= b \cdot ((\omega_1^*)^2 + (\omega_2)^2 - (\omega_3)^2 - (\omega_4^*)^2) = 0 \\ u_\psi &= b \cdot ((\omega_1^*)^2 - (\omega_2)^2 - (\omega_3)^2 + (\omega_4^*)^2) > 0 \\ u_f &= b \cdot ((\omega_1^*)^2 + (\omega_2)^2 + (\omega_3)^2 + (\omega_4^*)^2) > 0\end{aligned}$$

A counter clockwise yaw rotation can be executed by increasing  $\omega_2, \omega_3$ , then  $u_\psi > 0$ . The thrust force  $F$  has only a  $z$  component. A yaw rotational movement does not result in a translational movement.



**Fig. 2.3:** Sketch of an autonomous quadcopter flight control system

### 2.1.2 Autonomous Quadcopter

Autonomous quadcopter are especially interesting research topics, because it helps automating a lot of purposes. The general task can be in general defined as calculating the optimal motor signals  $S = \{s_1, s_2, s_3, s_4\}$  based on the mission goals  $g$  and the current pose  $p = \{p^*, a^*\}$  in order to achieve a stable flight and satisfy the goals. Therefore, it consists of an *inner* and an *outer control loop*.

The inner loop controls the attitude based on the current desired attitude  $a$  and the estimated attitude  $a^*$ . By subtracting these the attitude error  $e_a$  can be used to calculate the signals  $S$ , which are most commonly pwm signals. The signals are then applied by the motors, which act as actuators and apply a rotational force to the propellers. As a consequence the propeller get the rotational speeds  $\omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$  which directly influence the position of the drone in the environment. The state of the drone is captured by the sensor and the corresponding sensor values  $\sigma$ , which are falsified by noise and offsets. They are then used to estimate the attitude.

The outer loop controls mission goals based on the prior defined goals  $g$  and the current estimated position  $p$ . These inputs are used in order to calculate the desired attitude  $a$ .

In general Control loops in software are always time discrete. For cascaded controllers like here the inner loop is running at a higher frequency than the outer loop. This is crucial, because attitude control is very time sensitive and has to be implemented in real time. A bad or slow attitude control could lead to an unstable flight or in worst case it could crash and therefore automatically erase the possibility to achieve mission goals.

## Attitude Controller & Mixing

Common attitude controller often use *PID controller* as one of the most essential parts of nearly all commercial quadcopter flight control. PID controller are linear feedback controller with a proportional part, an integral part and a differential part and can mathematically described in two different ways:

$$u(t) = K_p \cdot e(t) + K_i \cdot \int_0^t e(\tau) d\tau + K_d \cdot \frac{de(t)}{dt} \quad (2.8)$$

$$u(t)^* = K \cdot (e(t) + \frac{1}{T_n} \cdot \int_0^t e(\tau) d\tau + T_v \cdot \frac{de(t)}{dt}) \quad (2.9)$$

However, since both are mathematically the same just with other gain definitions Equation (2.8) will be the relevant one hereafter. The control signal  $u(t)$  is calculated with the use of the configurable constant gains  $K_p, K_i, K_d$ , which weigh the different parts of the controller.

Within a PID controller the proportional term considers the current error  $e(t)$ , the integral term considers the history of errors by integrating about them and the differential term estimates the future error by considering change of the error. Since the attitude control is time discrete in software, Equation (2.8) must be transformed with the use of the sampling time period  $T$ :

$$u(t) = K_p \cdot e(t) + K_i \cdot T \cdot \sum_0^t e(\tau) + K_d \cdot \frac{e(t) - e(t-1)}{T} \quad (2.10)$$

In a quadcopter there is a PID controller for each of the three axis(roll, pitch, yaw), which are all working parallel during an inner loop control cycle. Then mixing is used to translate the PID values of each axis to a pwm signal  $s_i$  for each motor. This process uses a table, which consists of constants that describe the geometry of the frame. The throttle coefficient  $f$  and the mixer values of the motor  $M_i$  ( $m_{i,\phi}, m_{i,\Theta}, m_{i,\psi}$ ).

$$s_i = f \cdot (m_{i,\phi} \cdot u_\phi + m_{i,\Theta} \cdot u_\Theta + m_{i,\psi} \cdot u_\psi) \quad (2.11)$$

This implementation is modern state of the art attitude control. It is easy to implement, because Equation (2.10) and Equation (2.11) are simple equations and Equation (2.11) is just in need of some constants. In addition, there are not many calculations, so the process can be implemented in realtime. Also, this classic approach shows close to ideal performance in stable environments.

However, [KMWB19] shows that with RL an intelligent agent can be trained with the use of PPO (Section 2.2.5) that outperforms a classic PID Attitude controller in harsh, unpredictable environments.



## Position Estimation

Position or state estimation can be accomplished in different ways. By the use of classic filters the sensor values are more accurately, although there still is a relevant error. *Kalman Filters* are the common state of the art solution for position estimation, because they are statistical optimal. Although the name suggests otherwise a Kalman Filter is not a classic filter, but a stochastic weighted combination of the state prediction and measurement.

At each step, the next position and the next measurement is predicted. Based on the prediction of the measurement and the real measurement an innovation is calculated and used in combination with the Kalman Filter Gain  $K$  to calculate the next state. At the same time, the state has a covariance and the covariance of the next state is calculated. Then the covariance of the innovation and the Kalman Filter Gain  $K$  is calculated and used to correct the covariance of the state.

The amount of calculations increases exponentially with the amount of used sensor values. As a consequence, the use of sensors for the Kalman Filter is limited in realtime tasks. In addition it is in need of linearization. As a consequence there is research to further enhance the filter. For example, [SZWJ22] shows Deep Kalman Filters, which combine the learning ability of deep learning method and the noise filtering ability of the classic Kalman Filter can further enhance trajectory estimation. Although it should be mentioned that the mentioned work uses external satellite images of moving targets. So, the use case is quite different, but in theory adaptable.

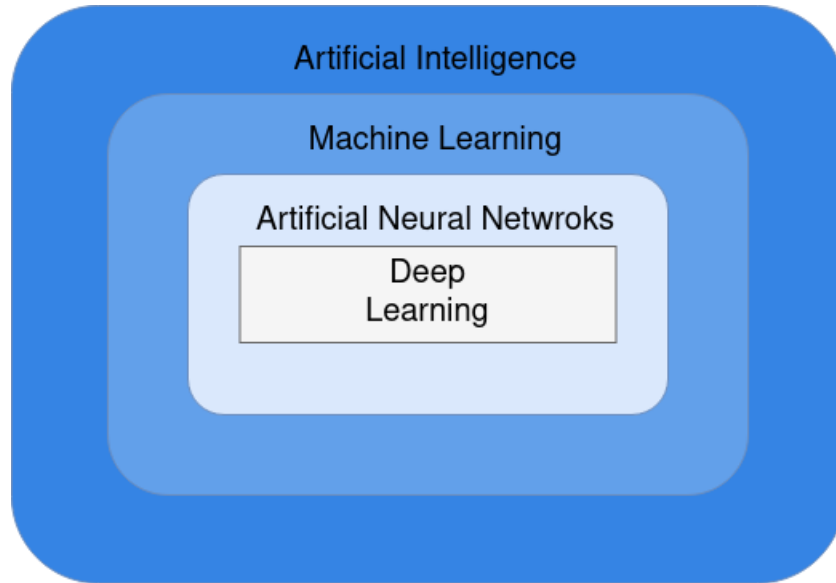
## Pathfollowing Controller

Pathfollowing Controller work in the outer control loop as Mission Goal Controller (Figure 2.3) with the goal  $g$  typically defined as a list of waypoints  $\Lambda = \{\lambda_0 \dots \lambda_i \quad i \in [0 \dots n]\}$ . These waypoints can be defined in different ways:

- $\lambda_i = \{x, y, z\}$ : classic point in 3D space
- $\lambda_i = \{x, y, z, \Theta, \phi, \psi\}$ : point in 3D space with attitude
- $\lambda_i = \{x, y, z, t\}$ : classic point in 3D space with time when it should be reached  
→ would lead to a Trajectoryfollowing Controller
- $\lambda_i = \{x, y, z, \Theta, \phi, \psi, t\}$ : point in 3D space with attitude and time  
→ would lead to a Trajectoryfollowing Controller

At the moment there is a wide range of different approaches for this task. Typically, either straight line paths or circular orbit paths are implemented, depending on the use case.

For example, [NBMB07] implemented a *Vector Field Pathfollowing* with promising results, although [SSS13] shows that this accuracy is combined with large control effort. But there is also the possibility to use a PID like implementation or a classic carrot chasing algorithm.



**Fig. 2.4:** Concept onion of Artificial Intelligence [DJS20]

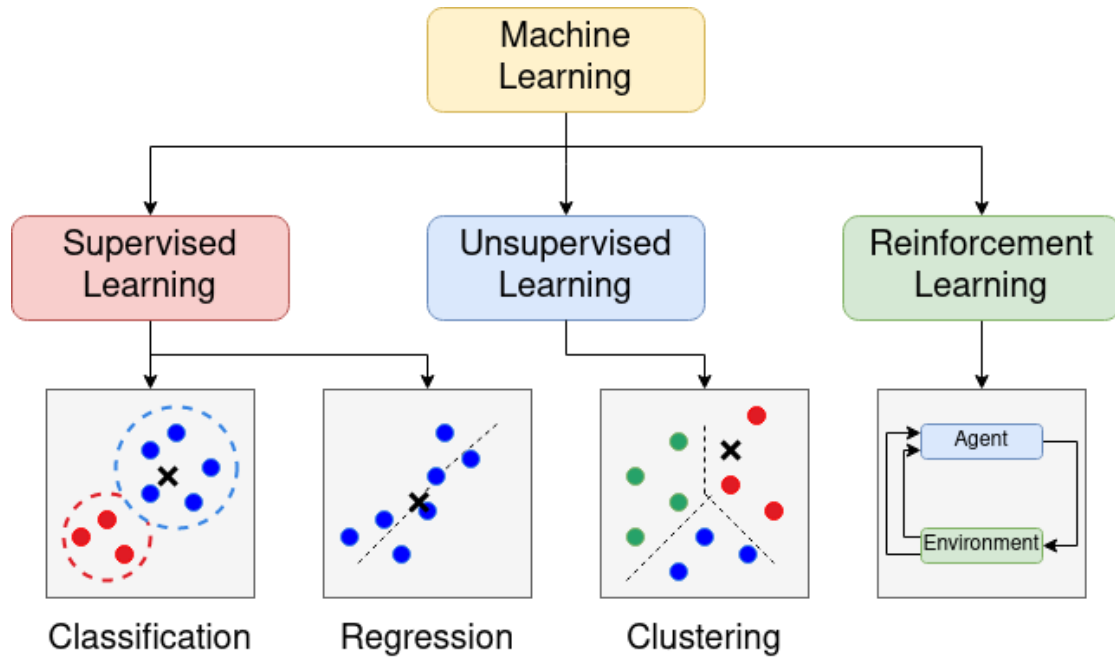
## 2.2 Reinforcement Learning Basics

In modern society the field of artificial intelligence is increasingly recognized and used. Be it transformer technologies such as ChatGPT or classic deep learning, artificial intelligence has applications in a variety of different fields of science and everyday life. Therefore, before explaining RL the different terms and aspects of the different methods should be distinguished.

*Artificial Intelligence* refers to an area in computer science, which deals with the simulation and application of behaviour that humans would define as intelligent. Although it is hard to define intelligence, there are some easy recognizable symptoms such as solving problems independently or human like interactions. The *Turing Test* [Tur50] represents an operative definition of intelligence. To sum it up, an algorithm passes the test, if a human, who previously asked written questions, can not determine whether the answer is given by the algorithm or a human. In order to pass this test, the algorithm is in need of different things:

- *Natural Language Processing (NLP)*: processing, interpretation, generation and output of natural language
- *Knowledge Representation*: understanding the content of the question
- *Logical Close*
- *Machine Learning (ML)*: adapt to new context, recognize structures

With these definitions and Figure 2.4 it is easy to locate RL, which is an area of machine learning and often uses artificial neural networks in order to satisfy the designated task.



**Fig. 2.5:** The main learning strategies of machine learning including the related approaches like classification and regression under the Supervised Learning, clustering under the Unsupervised Learning and a sketch of the Reinforcement Learning concept

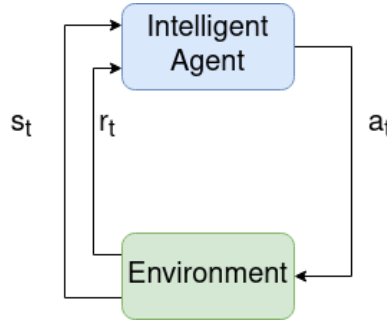
### 2.2.1 Classification of Machine Learning Areas

In general Machine Learning is seen as a subfield of artificial intelligence (Figure 2.4), that deals with generating knowledge based on experience. ML programs learn patterns based on training data and then can perform tasks without being explicitly programmed to do so by generalizing unknown examples.

#### Supervised Learning

*Supervised Learning* builds a model based on labelled data, which contains the input and the matching output. By iterative optimization supervised learning algorithms learn a function that can be used to predict the output associated with inputs. If learned correctly it also presents the output for inputs that were not a part of the training data set. In this case, the outputs should be analogue to the outputs of similar training data. The two most common types of Supervised Learning algorithms are *Classification* and *Regression* (Figure 2.5).

Classification algorithms are in need of outputs that are restricted to a limited set of categories. After training the model can determine which category belongs to a certain input. Regression algorithms learn a model to estimate the relationship between a dependant and one or more independent variable and are used when the outputs may have a numerical value within a defined range.



**Fig. 2.6:** The basic concept of RL: the agent takes an action based on the observed state and receives a numerical reward at each time step

## Unsupervised Learning

*Unsupervised Learning* builds a model with unlabelled data, which only contains the input data. In contrary to Supervised Learning the data set does not contain a preferred output for the given inputs. Unsupervised Learning aims at discovering a distribution of the data in order to gain new knowledge.

*Clustering* (Figure 2.5) discovers clusters within the given data and can therefore predict for new inputs the related cluster.

## Reinforcement Learning

*Reinforcement Learning* is an area of ML designated to modelling decision making by finding a way to take an action in an environment in order to maximize the sum of rewards. It has proven successful in a wide range of different problems like game theory, control theory or swarm intelligence. The main difference of RL to the already mentioned areas of ML is that the action affects the environment and therefore the observation it receives. Also, it learns by the use of a reward signal.

The *intelligent agent* is the decision maker and typically implemented with the use of a *Neural Network* (Section 2.2.4). Figure 2.6 shows the basic idea: at each time step  $t$  the agent takes an action  $a_t$  based on the prior observed state  $s_{t-1}$ . Based on how good the action was he receives a reward  $r_{t+1}$  and the new state  $s_{t+1}$ . The reward signal can either depend on the state ( $s_t$ ), the state and the chosen action ( $s_t, a_t$ ) or the state, the chosen action and the transitioned state ( $s_t, a_t, s_{t+1}$ ). It is always numerical and used to optimize the agent. The environment  $\epsilon$  receives at each time step the action  $a_t$ . Then, it simulates the state transition based on its implementation. It calculates the reward signal  $r_{t+1}$  and returns it alongside with the observed transitioned state  $s_{t+1}$ . It should be mentioned that the observed state  $s_t$  must not be the complete internal state of the environment  $s_t^*$ . Furthermore, it is only an observation and also can be defined as  $o_t$ . In general, the environment can be modelled with the use of a *Markov Decision Process*.

### 2.2.2 Markov Decision Process

A *Markov Decision Process* (MDP) is a time discrete stochastic control process, which constitutes a mathematical framework for modelling decision making. A MDP is formally defined as a 5-Tupel  $(S, A, R, P, p_0)$ , where:

- $S$  is the state space with the states  $s_i \in S$
- $A$  is the action space with the actions  $a_i \in A$
- $R$  is the reward function, mapping  $S \times A \times S \rightarrow \mathbb{R}$  with:  
 $r_{t+1} = R(s_t)$ ,  $r_{t+1} = R(s_t, a_t)$  or  $r_{t+1} = R(s_t, a_t, s_{t+1})$
- $P$  is the state transition probability function mapping  $T(s_t, a_t, a_{t+1}) \sim Pr(s_{t+1}|s_t, a_t)$
- $p_0$  is the starting state distribution

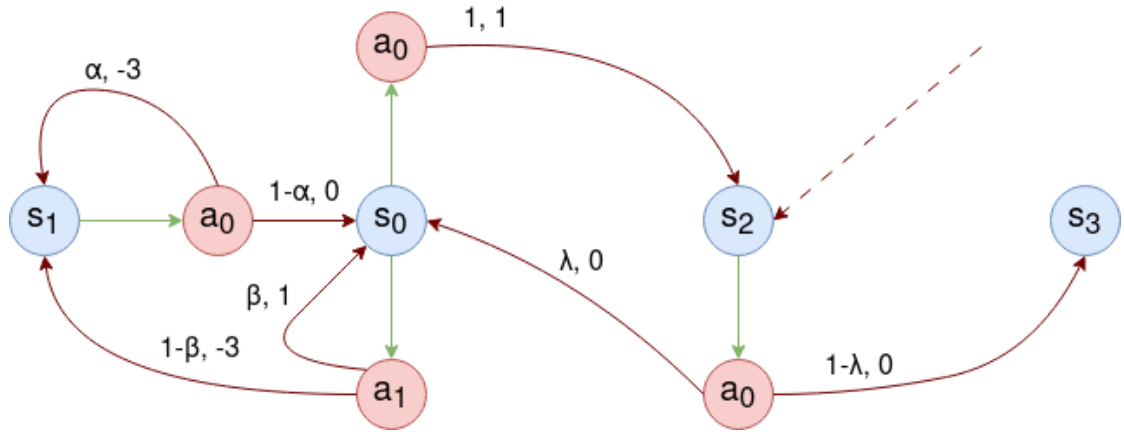
The state space  $S$  is a set of states  $s$  and may be discrete or continuous. Analogue, the action space  $A$  may be discrete or continuous. If there are final states the MDP is called *finite MDP*, else *continuing MDP*.

The reward function always returns a numerical reward. Like previously mentioned the reward may only depend on the state  $s_t$  or state and action  $a_t$ , or state, action and transitioned state  $s_{t+1}$ .

When transitioning from one state to the next, the transitioning must not be deterministic.  $P$  maps a probability to a transition  $T(s_t, a_t, a_{t+1})$  from state  $s_t$  to state  $s_{t+1}$  by choosing the action  $a_t$ . But there also may be deterministic transitions. Then there only is one possible state to transition to and taking the related action will result in this transition with a probability of 1. In addition, the transitions always support the Markov property, that defines that transitions only depend on the most recent state and action. So the prior states  $s_0 \dots s_{t-1}$  are irrelevant for the current transition.

In addition, the starting state distribution defines which states may be the first state  $s_0$ .

By modelling a RL problem as a MDP, there is a mathematical basis to optimize the decision making. Further, it defines the possible states and actions. Choosing a good reward function is the key to influence the decision making. For example, by penalizing the transition to a bad state the agent will learn to avoid to choose the matching action. Analogue, a big positive reward can influence the decision making.



**Fig. 2.7:** Visualization of a simple MDP with 4 states ( $s_0, s_1, s_2, s_3$ ) (blue nodes), 2 actions ( $a_0, a_1$ ), state-action edges (green) and action-state edges (dark red) containing a tuple (possibility, reward)

### Visualization

A Markov Decision Process can be visualized with the help of a directed transition graph. The graph possesses state nodes and action nodes. Each state node is connected with  $n \in \mathbb{N}_0$  actions with directed state-action edges. Each action node is connected with  $k \in \mathbb{N}$  states with action-state edges, that possess a probability  $\alpha_i$ . Since after taking in action there must be a transition, it follows that the sum of all the outgoing action-state edge probabilities must be equal to 1:

$$\sum_{i=0}^k \alpha_i = 1 \quad (2.12)$$

It is also possible that there is a state-action edge from state node  $s_k$  to action node  $a_n$  and an action-state edge from  $a_n$  back to  $s_k$ . It is equivalent to taking an action that may cause no transition at all and the state stays the same. In addition the graph may pose starting and ending states.

For example, Figure 2.7 shows a finite MDP with 4 states (blue nodes).  $s_2$  is defined as starting state and since no action can be taken in state  $s_3$ , this is an ending state. If a state possesses multiple state-action edges like  $s_0$  a decision has to be made. After choosing an action there may be multiple action state edges like action  $a_1$  from state  $s_0$ . In this example there is a state transition to  $s_1$  with a probability of  $1 - \beta$  that succeeds in a numerical reward of  $-3$ . Also, there is the possibility of  $\beta$  to transition into  $s_0$  which causes a reward of 1.

### 2.2.3 V,Q Function & The Bellman Equations

In order to further define how the agent learns, there are a couple of functions, equations and helpful definitions. The *policy*  $\pi$  defines which action should be taken based on the current state and the *trajectory*  $\tau$  defines different possible sequence of transitions inside the MDP. In order to find the optimal trajectory there are different metrics defining an expected value to each state (*V Function* Section 2.2.3) and each state-action pair (*Q Function* Section 2.2.3). With these functions different policies can be evaluated.

#### Policy

The policy  $\pi$  is a mapping from the state space  $S$  to the action space  $A$ . Formally this mapping is described as

$$S \rightarrow \mathbb{P}(A, \pi(a_t, s_t))$$

and defines the probability of choosing the action  $a_t$  in state  $s_t$ .

In many RL Algorithms *Neural Networks* (Section 2.2.4) will be used to model the policy with the goal to approximate an optimal policy  $\pi^*$  that maximizes the expected reward:

$$\pi^* = \arg \max_{\pi} J(\pi)$$

#### Trajectory

A sequence of states  $s_0 \dots s_{T+1}$  and the matching action  $a_0 \dots a_T$  is called a trajectory  $\pi$  with the amount of steps  $T$ . For example there could be the trajectory  $\tau_1 = (s_2, a_0, s_0, a_1, s_1)$  for the visualized MDP (Figure 2.7) with  $T = 2$  steps: Since  $s_2$  is defined as a starting state it always is the first state of the trajectory. By choosing  $a_0$  there is a state transition to state  $s_0$ . Then there is a transition to state  $s_1$  by choosing the action  $a_1$ . This trajectory  $\tau_1$  could be achieved with the policy  $\pi_1$  with  $\pi_1(a_0|s_2) = 1, \pi_1(a_1|s_0) = 1$ . This policy literally defines that in the state  $s_2$  always the action  $a_0$  should be taken and in state  $s_0$  always the action  $a_1$ . Further there can be described a probability that this trajectory  $\tau$  occurs under the policy  $\pi$ .

$$P(\tau|\pi) = p_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi(a_t|s_t) \quad (2.13)$$

This is basically the the product of the probabilities that the policy  $\pi$  takes an action and that the action  $a_t$  transitions the state to  $s_{t+1}$ . In the above example, this would result in

$$P(\tau_1|\pi_1) = 1 \cdot \lambda \cdot 1 \cdot (1 - \beta) = \lambda \cdot (1 - \beta)$$

**Tab. 2.1:** Example of all two step trajectories with the matching probability, reward and expected reward of the MDP shown in Figure 2.7

| $\tau$                      | $P(\tau \pi)$               | $R(\tau)$ | $J(\tau)$                            |
|-----------------------------|-----------------------------|-----------|--------------------------------------|
| $(s_2, a_0, s_3)$           | $(1 - \lambda)$             | 0         | 0                                    |
| $(s_2, a_0, s_0, a_1, s_0)$ | $\lambda \cdot \beta$       | 1         | $\lambda \cdot \beta$                |
| $(s_2, a_0, s_0, a_1, s_1)$ | $\lambda \cdot (1 - \beta)$ | -3        | $-3 \cdot \lambda \cdot (1 - \beta)$ |

### Discounted & Undiscounted Reward

In *finite*, episodic MDPs with ending states, the sum of the rewards is easy to calculate for a trajectory  $\tau$  with the use of the amount of steps  $T$ .

$$R(\tau) = \sum_{t=0}^{T-1} r_{t+1} \quad (2.14)$$

It was already discussed in Section 2.2.2 that a MDP may be infinite and continuing. As consequence the undiscounted reward, that was defined in the above may diverge, because there is no ending state. In order to counter the diverging of the sum a discount factor  $\gamma \in [0...1]$  is used. The idea behind is to weigh future rewards. With increasing future step  $t$  the weigh factor decreases exponentially. As a consequence the sum of the rewards converges and is called discounted reward.

$$R(\tau)_t = \sum_{k=0}^{\infty} \gamma^k \cdot r_{t+k+1} \quad (2.15)$$

$$= r_{t+1} + \gamma \cdot r_{t+2} + \gamma^2 \cdot r_{t+3} + \dots \quad (2.16)$$

Equation (2.15) shows how the discounted reward is calculated:

Future rewards are weighed with the factor  $\gamma^n$   $n \in \mathbb{N}_0$  depending on how many steps they are in the future. The special case of  $\gamma = 1$  means that the task should be episodic, the case of  $\gamma = 0$  results in a greedy implementation, where future rewards are not considered at all and the agent learns to always choose the action that implies the highest reward in that step. With the before mentioned possibility of a transition  $\tau$  under the policy  $\pi$  an expected reward  $J(\pi)$  of a policy can be defined. For example in Figure 2.7 there can be a couple of  $T = 2$  step trajectories  $\tau_i$  with a matching probability and reward  $R(\tau_i)$  (Table 2.1) that leads to an expected reward  $J(\pi_1)$  of depending on the transition probabilities  $\lambda, \beta, \alpha$ :

$$J(\pi_1) = \lambda \cdot \beta - 3 \cdot \lambda \cdot (1 - \beta)$$



## V-Function

A vast majority of *Reinforcement Learning Algorithms* try to estimate a *Value Function* in order measure how good it is to be in a state  $s_t$  under a matching policy  $\pi$ . The V Function calculates the expected return  $J(\pi)$  with the starting state  $s_t$  and acting accordingly to the defined policy:

$$V^\pi(s) = \mathbb{E}_{\tau \sim P(\cdot|\pi)}[R(\tau)|s_0 = s] \quad (2.17)$$

$$= \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s=s_t\right] \quad (2.18)$$

In addition, an *optimal Value Function* (Equation (2.19)) can be defined and used in order to calculate an optimal policy:

$$V^*(s) = \max_{\pi} V^\pi(s) = \max_{\pi} \mathbb{E}_{\tau \sim P(\cdot|\pi)}[R(\tau)|s_0 = s] \quad (2.19)$$

## Q-Function

The *Q Function* (action-value function) is defined quite similar to the Value Function but instead of starting in a state, the Q Function starts in a state-action pair and then takes actions according to the chosen policy. Therefore, it is also mathematically defined as the expected reward:

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim P(\cdot|\pi)}[R(\tau)|s_0 = s, a_0 = a] \quad (2.20)$$

Analogue to the V Function, an *optimal Q Function* can be calculated by taking the maximal expected reward.

$$Q^*(s, a) = \max_{\pi} Q^\pi(s) = \max_{\pi} \mathbb{E}_{\tau \sim P(\cdot|\pi)}[R(\tau)|s_0 = s, a_0 = a] \quad (2.21)$$

With the help of the optimal Q Function the optimal action  $a$  for every state can be chosen by comparing their optimal Q values and used to construct an optimal policy  $\pi^*$ . If the MDP has discrete actions and states this is no problem and then called *Tabular Reinforcement Learning*. For more complex MDPs RL algorithms that use NN can be used. Probably the best known algorithm of this kind is *Deep Q Learning* [FWXY20] which estimates the Q Function with the use of a behaviour and a target policy. These policies are implemented with the use of Neural Networks.

## The Bellman Equations

The *Bellman equations* are the central point of some RL algorithms. It parts the V and Q function into two separate parts consisting of the instant reward and the discounted future values. As an effect it simplifies the computation of the value function. Instead of building a sum over multiple time steps, an optimal solution of a complex problem is found by finding optimal solutions for simpler, recursive subproblems.

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s), s' \sim P(\cdot|s,a)}[R(s, a, s') + \gamma \cdot V^\pi(s')] \quad (2.22)$$

$$Q^\pi(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)}[R(s, a, s') + \gamma \cdot \mathbb{E}_{a' \sim \pi(\cdot|s')}[Q^\pi(s', a')]] \quad (2.23)$$

Equation (2.22) shows the instant reward  $R(s, a, s')$  which is computed by calculating the expected value for all possible actions  $a$  and all possible transition states  $s'$  and the discounted future  $\gamma \cdot V^\pi(s')$  reward of the values of the possible transition states. Analogue, Equation (2.23) shows the split up of instant reward and future reward for state-action pairs. As a consequence, iterative approaches can be implemented, that possesses the ability to calculate the values of all states or the values of all state-action pairs.

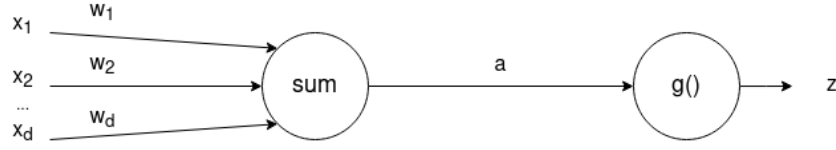
## The Bellman Equations of Optimality

It was already discussed how optimal V and Q functions can be calculated and used to determine the optimal actions and create an optimal policy (Equation (2.19), Equation (2.21)). Analogue the *Bellman Equations of Optimality* can be defined with the use of Equation (2.22) and Equation (2.23):

$$V^*(s) = \max_a \mathbb{E}_{s' \sim P(\cdot|s,a)}[R(s, a, s') + \gamma \cdot V^*(s')] \quad (2.24)$$

$$Q^*(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)}[R(s, a, s') + \gamma \cdot \max_{a'} Q^*(s', a')] \quad (2.25)$$

## Monte-Carlo & TD-Sampling



**Fig. 2.8:** The McCulloch-Pitts model of a single neuron uses a weighted sum with the inputs  $x_i$  and weights  $w_i$  and a non-linear activation function  $g()$  in order to compute the output  $z$  [Bis94]

## 2.2.4 Neural Networks

A *Neural Network (NN)* can be described as an reactive computing system that consists of multiple simple *Neurons* that are interconnected. Therefore, the output does depend on the inputs and the defined dynamic state response of each neuron. NN are inspired by the network of neurons in the human brain which pass information to another with the help of synapses. While biological brains are dynamic and analog, NN tend to be static and symbolic.

NN have shown to be very useful, because they are universal function approximators. In addition, they are adaptive and poses the ability to change their structure based on external or internal information that flows threw it. They are non-linear and can be used to model complex relationships between inputs and outputs or find specific pattern in data.

### Neuron

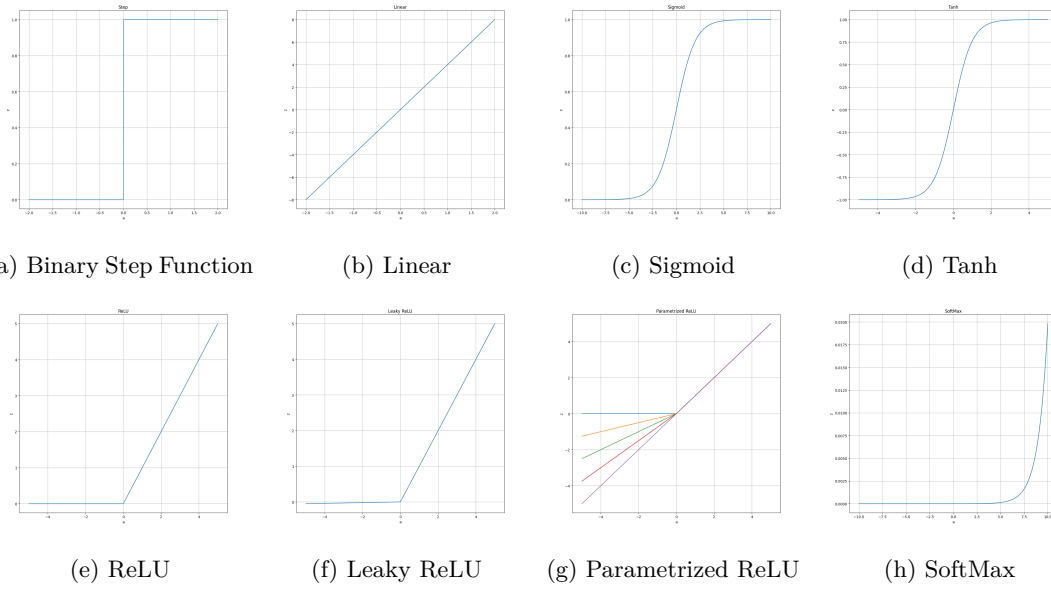
A Neuron is the basic element of a NN. Mathematically a neuron can be seen as a non linear function, which computes the output  $z$  based on a set of inputs  $x_i, i \in [1, \dots, d]$  [Bis94]. Further, [MP43] introduced the simple mathematical framework, called the McCulloch-Pits model (Figure 2.8):

$$a = \sum_{i=1}^d w_i \cdot x_i + w_0 \quad (2.26)$$

$$a = \sum_{i=0}^d w_i \cdot x_i \quad (2.27)$$

$$z = g(a) \quad (2.28)$$

Each input  $x_i$  is multiplied with its matching weight  $w_i$ , which is a parallel to synaptic strength in biological networks like brains (Equation (2.26)). In addition, there is a offset which is called *bias* and is analogue to the firing threshold of a biological neuron. By using the bias as additional input and setting it to 1, the sum can be further simplified (Equation (2.27)). The output of the neuron is then computed by giving it to a non-linear *Activation Function*  $g()$ .



**Fig. 2.9:** 8 different Activation Functions.

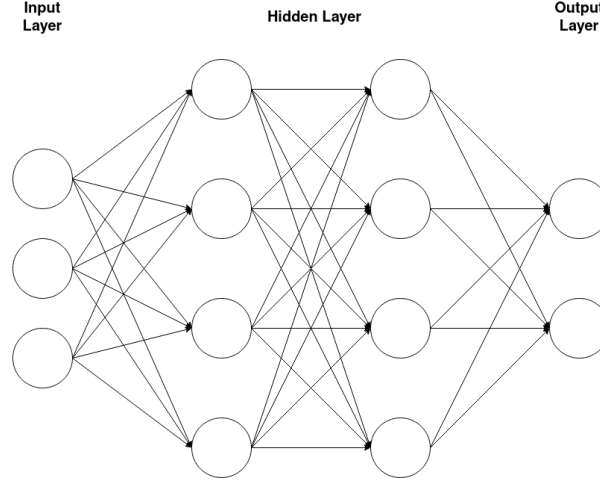
## Activation Function

It was already discussed that NN are universal function approximators and that the neuron consists of a linear sum (Equation (2.27)). In order to represent nonlinear convoluted functional mappings between input and output, non linearity is added to each neuron with a *Activation Function*. It is important that these functions are differential in order to implement a back propagation optimization strategy [SSA17].

At the moment, there are different activation functions commonly used (Figure 2.9):

- Binary Step Function
- Linear
- Sigmoid
- Tanh
- ReLU
- Leaky ReLU
- Parametrized ReLU
- SoftMax

Each of this different functions has its advantages. For example, SoftMax is good at multiclass classification problems, while Sigmoid is good at binary classification problems. Also, in a NN with ReLU not all neurons are activated at the same time, which increases efficiency.



**Fig. 2.10:** Simple Feed-Forward Network with 3 neurons in the input layer, 2 hidden layers with 4 neurons and with 2 layers in the output layer

### Simple Feed-Forward Network

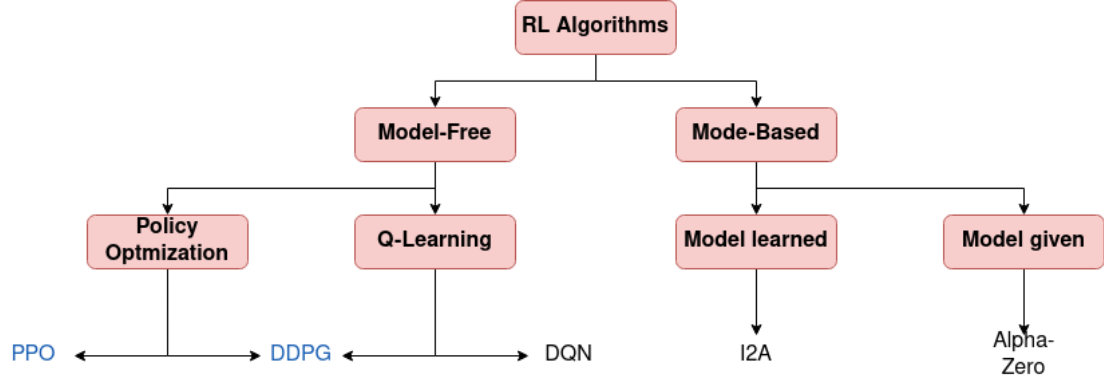
A *Feed-Forward Network* is an artificial network wherein information directly flows forward in one direction from the input to the output. In contrast to *Recurrent Neural Networks*, it does not contain cycles of nodes and information can not flow in loops. It can be described as a directed graph  $G = (V, E, w)$  with neurons as the set of nodes  $V$  and the set of edges  $E$  with the matching weight  $w$  (Figure 2.10). The neurons are grouped into different Layers:

- **Input Layer:** each node receives its input  $x$  and calculates the output  $z$  with the use of Equation (2.27) and passes it to each node of the next layer
- **Hidden Layer:** each node receives its inputs  $x_i, i \in [1...d]$  from the  $d$  nodes of the previous layer and uses Equation (2.27) and the weights  $w_i, i \in [1..d]$  in order to calculate the output  $z$  and pass it to each node of the next layer
- **Output Layer:** each node receives its inputs  $x_i, i \in [1...d]$  from the  $d$  nodes of the previous layer and uses Equation (2.27) and the weights  $w_i, i \in [1..d]$  in order to calculate the output  $z$  of the NN

The number of hidden layers can differ, likewise the amount of neurons in each layer. Due to the structure of neurons even simple Feed-Forward Networks can approximate functions  $F : \mathbb{R}^j \rightarrow \mathbb{R}^k$ , where the dimension  $j$  equals the amount of nodes in the input layer and  $k$  equals the amount of nodes in the output layer.

### Deep Neural Network

*Deep Neural Networks* differ from the previous mentioned NN in terms of complexity. Due to the large amount of neurons and a lot of hidden layers, the NN is able to approximate very complex mappings from the input to the output space.



**Fig. 2.11:** Overview of the different categories of RL algorithms with examples: PPO and DDPG Algorithm coloured because they were used in this work.

### 2.2.5 State of the Art Algorithms

Modern state of the art RL algorithms can be divided into different categories (Figure 2.11):

*Model-based Algorithms* uses its experience in order to construct a model of the environment by modelling the MDP transitions. With the use of this model the optimal actions are chosen for the policy. This category can be further divided into *model learning algorithms* like I2A and *model given algorithms* like Alpha-Zero. Model learning algorithms observe the trajectory  $\tau$  while following a policy  $\pi$  and use it in order to learn a dynamics model. Then the algorithms plan through the dynamics model in order to choose the actions. Model given algorithms work quite similar but the dynamic model is already given.

In *Model-free Algorithms* the agent constructs a policy based on trial-and-error experience and can be further divided into *Policy Optimization* and *Q-Learning*. Policy Optimization algorithms like *PPO* the policy is learned directly. In Contrast Q-Learning algorithms like *Deep Q-Learning* learns the Q-function and updates it in order to get an optimal Q-function. In addition there are hybrid algorithms like *DDPG* that combines learning the policy function and Q-function.

RL algorithms can be further categorized into *on-policy* and *off-policy* algorithms. On-policy algorithms use the current optimized policy in order to gain experience that is then used to further optimize the policy. In contrast, off-policy algorithms learn from actions that may not be according to the current policy. This experience is then used to construct the optimal policy.

The algorithms used in this work are part of *Stable Baselines 3* [RHG<sup>+</sup>21], which is a set of reliable implementations of RL algorithms in *PyTorch* that promise efficient learning and a good base to build projects on top of. Each of the main state of the art RL algorithms can be found there and used with custom policies, custom environments and custom callback functions. In addition, it offers *Tensorboard* support, which is a tool that can be used to visualize the training process.

## PPO - Basics

PPO algorithms [SWD<sup>+</sup>17] are a family of policy gradient methods for RL, they directly try to improve the policy  $\pi$  as much as possible in a single step without stepping to far and avoiding performance collapse. It is an *Actor-Critic Method* that approximates the value function  $V_\phi$  and the policy  $\pi_\Theta$  with the matching set of parameters  $\phi, \Theta$ . In each step the actor updates the policy parameters and the critic updates the value function parameters. In order to measure how good a policy  $\pi_\Theta$  performs in relation to an old policy  $\pi_{\Theta_{old}}$  a surrogate objective is used, that keeps the new policies close to the old one:

$$L^{PPO}(\Theta) = \mathbb{E}_{\tau \sim P(\cdot | \pi_{\Theta_{old}})} \left[ \sum_{t=0}^T L(s_t, a_t, \Theta_{old}, \Theta) \right] \quad (2.29)$$

$$L(s_t, a_t, \Theta_{old}, \Theta) = \min \left( \frac{\pi_\Theta(a_t | s_t)}{\pi_{\Theta_k}(a_t | s_t)} A(s_t, a_t), g(\epsilon, A(s_t, a_t)) \right) \quad (2.30)$$

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon) \cdot A & \text{if } A \geq 0 \\ (1 - \epsilon) \cdot A & \text{if } A < 0 \end{cases} \quad (2.31)$$

This surrogate objective mainly depends on the hyperparameter  $\epsilon$  which is an upper bound to the distance between the policies and the advantage function  $A$ , that estimates how good an action is compared to the average action for a specific state.

Equation (2.29) can be explained intuitively by observing a single state-action pair:

If the advantage of the state-action pair is positive Equation (2.30) is reduced to

$$L(s_t, a_t, \Theta_{old}, \Theta) = \min \left( \frac{\pi_\Theta(a_t | s_t)}{\pi_{\Theta_k}(a_t | s_t)}, (1 + \epsilon) \cdot A(s, a) \right) \quad (2.32)$$

Because of the positive advantage the objective will increase with  $\pi_\Theta(a|s)$ , but is still limited to  $(1 + \epsilon) \cdot A(s, a)$  by the minimum expression.

If the advantage of the state action pair is positive Equation (2.30) is reduced to

$$L(s_t, a_t, \Theta_{old}, \Theta) = \max \left( \frac{\pi_\Theta(a_t | s_t)}{\pi_{\Theta_k}(a_t | s_t)}, (1 - \epsilon) \cdot A(s, a) \right) \quad (2.33)$$

Because of the negative advantage the objective will increase with the decrease of  $\pi_\Theta(a|s)$ , but is still limited to  $(1 - \epsilon) \cdot A(s, a)$  by the maximum expression.

As a consequence the policy does not profit from changing a lot.

PPO is an on-policy algorithm, so the actions are selected according to the latest policy. Therefore, the randomness of the action selection decreases over time, which may causes the problem of being trapped in a local optima.

## PPO - Algorithm

The PPO algorithm (Algorithm 1) first initializes the policy and value net with the matching parameters  $\Theta_0$  and  $\phi_0$ .

Then the following steps are then executed until convergence:

By running the current policy  $\pi_{\Theta_k}$  a set of Trajectories  $D_k$  is collected and the rewards  $R_t$  are computed. Then the advantage  $A$  is computed based in the current value function approximation  $V_{\phi_k}$ . With the use of this values the surrogate objective (Equation (2.29)) is calculated and used to update the policy  $\pi$ . Typically stochastic gradient ascent is used in order to update the policy. At last the value function is updated by using regression on the mean squared error.

Since the release of PPO it is known to be a state of the art algorithm with a good performance that is at least as good as other Policy Optimization methods on a variety of environments. In addition, the base algorithm is quite easy to implement.

Nevertheless there are still known problems:

The algorithm is quite sensitive to the initialization parameters  $\Theta_0, \phi_0$ . Like previously mentioned PPO can get stuck in local optima. This is the case if there are local optimal actions close to initialization.

PPO seems to be unstable when the reward function is not bounded on continuous action spaces.

---

### Algorithm 1: Proximal Policy Optimization [SWD<sup>+</sup>17]

---

```

1 Initialize policy parameters  $\Theta_0$  and value function parameters  $\phi_0$ 
2  $k = 0$ 
3 repeat
4   Collect set of trajectories  $D_k = \{\tau_i\}$  by running policy  $\pi_{\Theta_k}$ 
5   Compute rewards-to-go  $R_t$ 
6   Compute advantage estimates  $A$  based on the current value function  $V_{\phi_k}$ 
7   Update the policy by maximizing the PPO objective

      
$$\Theta_{k+1} = \arg \max_{\Theta} \left( \frac{1}{|D_k|} \sum_{\tau_i \in D_k} \sum_{t=0}^T \min \left( \frac{\pi_{\Theta}(a_t|s_t)}{\pi_{\Theta_k}(a_t|s_t)} A(s_t, a_t), g(\epsilon, A(s_t, a_t)) \right) \right)$$


      typically via stochastic gradient ascent
8   Fit value function by regression on mean-squared error

      
$$\phi_{k+1} = \arg \min_{\phi} \left( \frac{1}{|D_k|T} \sum_{\tau_i \in D_k} \sum_{t=0}^T (V_{\phi}(s_t) - R_t)^2 \right)$$

9    $k += 1$ 
10 until convergence

```

---



## 2.3 Pybullet Physics Simulator & Gym Pybullet Drones

*Pybullet* is a Python simulator that provides a fast and easy to use module for machine learning and robotics simulation. It provides dynamics simulation, inverse dynamics computation, forward and inverse kinematics, collision detection and ray intersection queries [CB21].

*Gym Pybullet Drones* [PZZ<sup>+</sup>21] is an *OpenAI Gym environment* based on the previously mentioned Pybullet physics simulator as back-end. It aims at multi-agent reinforcement learning for quadcopter simulations. Therefore, it provides a small variety of dronemodels as urdf files, for example based on the *Bitcraze’s Crazyflie 2.x nano-quadcopter* and gym environments.

### 2.3.1 BaseAviary Class

In Gym Pybullet Drones all aviary classes inherit from the *BaseAviary Class*. It provides a couple of useful methods, the most important ones named hereafter:

- `init`: initializes environment and loads the drone urdf file
- `reset`: resets the environment
- `step`: simulates one step and returns the observation, reward, a done information and an info
- `render`: renders the environment as textual output
- `startVideoRecording`: starts the recording of a mp4 video
- `physics`: the base PyBullet physics implementation
- `dynamics`: the drone dynamics implementation

### 2.3.2 BaseSingleAgentAviary Class

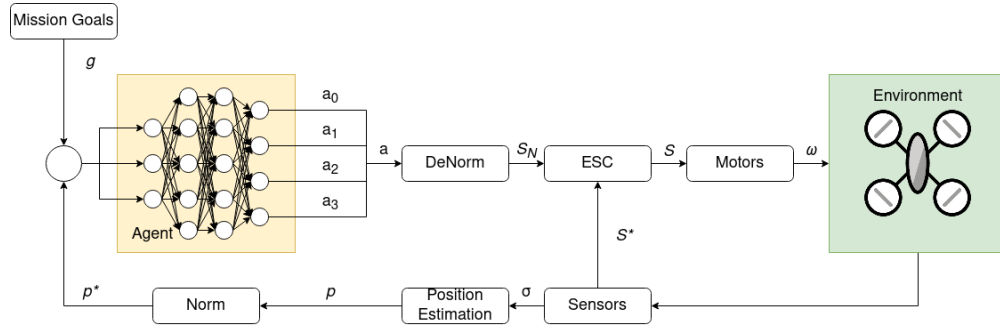
The *BaseSingleAgentAviary Class* is the most important environment class for this work, since it provides a framework for single agent RL problems. It inherits from the *BaseAviary Class* and defines typical action types, observation types and preprocesses the action that is passed to the `step` method of its mother class.

The implemented environment in this work inherits from this class.

### 3 Flight-Control Concept

It was already discussed in Section 2.1.2 how the flight control of modern UAV is structured. Since this work aims at using reinforcement learning to learn a robust control of a quadcopter, the explained does not fit this structure. Figure 2.3 shows the use of a mission goal controller that shows some distant resemblance to the implemented intelligent agent. In Contrast, it does not output a desired attitude to achieve the defined goal, but directly outputs actions that are more or less the motor values. As a consequence there is no need for a classic Attitude Control but a need of controlling the motor signals to match the actions. This chapter proposes a flight control concept with the use of an intelligent agent. Also, it provides a general concept of the implementation in order to get a wider view at the different implemented classes, tools and scripts.

At each timestep  $t$  the agent receives a mission goal  $g = (g_x, g_y, g_z)$  and a normalization of the position/state estimation  $p = (p_x, p_y, p_z, p_{roll}, p_{pitch}, p_{yaw}, \dot{p}_x, \dot{p}_y, \dot{p}_z, \dot{p}_{roll}, \dot{p}_{pitch}, \dot{p}_{yaw})$  and outputs a 4-tupel of actions  $a = (a_0, a_1, a_2, a_3)$  with each  $a_i \in [-1, 1], i \in [0, 1, 2, 3]$ . This can be denormalized to a signal  $S_N$ . On a real UAV this would be most likely a pwm signal. With the use of a speed controller the signal is held at the specific value. With the use of the motors and props there is a thrust  $f$  and a corresponding state transition that can be observed by the sensors. These sensor values  $\sigma$  can be used in order to estimate the position and normalize it to the range  $[-1, 1]$ . This concept is visualized in Figure 3.1.



**Fig. 3.1:** Proposed Concept of autonomous flight control with the use of an intelligent agent implemented with the use of a NN. The output of the NN is denormalized in order to translate them to motor signals. With the use of the sensors and position estimation the input for the next control step are calculated.

## 4 Implementation

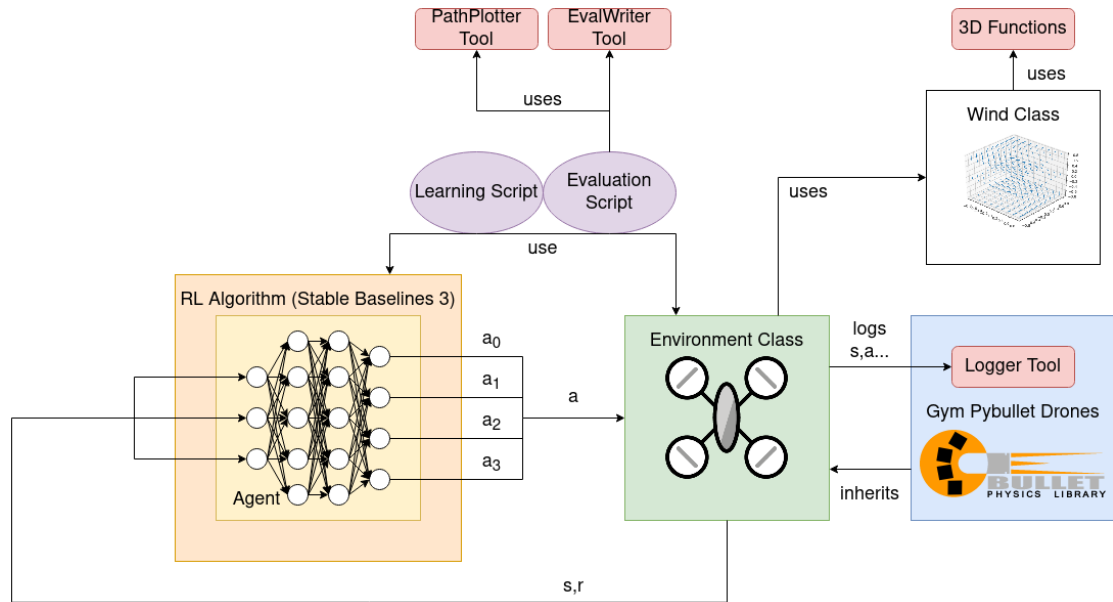
The concept of implementation can be divided into the used packages (*Stable Baselines3*, *Gym Pybullet Drones*, ...), the scripts, the environment classes and different support classes (Figure 4.1).

*Stable Baselines3* is used as implementation of the PPO algorithm (Algorithm 1) in order to achieve the tasks. *Gym Pybullet Drones* has already some environments and is the foundation of the simulation. It already provides a suitable step function and a well designed physics engine.

The scripts (Section 4.2 , lila) are used to either learn the agent on a given environment or evaluate it.

The environment classes(Section 4.1, green) inherit from a *Gym Environment* and models a MDP. By modelling this MDP precisely, it is defined what is learned later by the intelligent agent. It uses the wind class in order to model a harsh environment. In addition, there are a couple of evaluation tools (Section 4.2.3, red), that supports the implemented classes.

The whole concept is implemented in Python3.x with the use of a *conda environment*.



**Fig. 4.1:** Concept of the implemented software: the different tools(red), scripts(violet) that are used in order to learn the intelligent Agent robust flight control with the use of RL.

## **4.1 Environment Classes**

### **4.1.1 WindSingleAgentAviary Environment Class**

## Wind Class

## Modes

**State Space**

**Observation Space**

**Tab. 4.1:** The different ActionTypes with the corresponding dimensionality of the action, its range and how it is processed.

| ActionType       | dim       | range             | processing  |
|------------------|-----------|-------------------|---|
| <i>one_d_rpm</i> | $ a  = 1$ | $a_i \in [-1, 1]$ | $rpm = (hover\_rpm \cdot (1 + \alpha \cdot a)) \cdot (1, 1, 1, 1)$                      |
| <i>rpm</i>       | $ a  = 4$ | $a_i \in [-1, 1]$ | $rpm = hover\_rpm \cdot (1 + \alpha \cdot a)$   |
| <i>vel</i>       | $ a  = 4$ | $a_i \in [-1, 1]$ | $rpm = pid(S, vel = limit \cdot  a_3  \cdot \frac{(a_0, a_1, a_2)}{ (a_0, a_1, a_2) })$ |

## Action Space

The WindSingleAgentAviary Environment possesses three different type of action spaces, that are processed in different ways to the *rpm* (*rotation per minute*) of the four motors (Table 4.1). Nevertheless, all action types are continuous and are ranged in  $[-1, 1]$ .

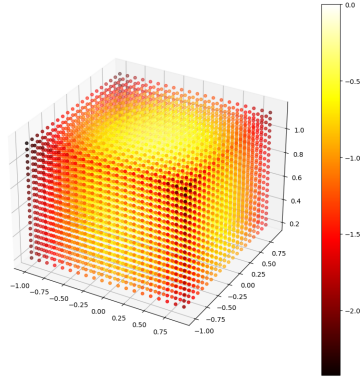
*one\_d\_rpm* is a one dimensional action space. The chosen action  $a$  is processed to range of  $\alpha$  about the *hover\_rpm* to a 4-tupel of rpms. The *hover\_rpm* is defined as the rpm that corresponds to hovering (Section 2.1.1). The 4-tupel is then forwarded to the motors. As a consequence of this limitation, the drone can only perform hovering, rising and falling movements and can not influence its  $x, y$  position or  $\Theta, \phi, \psi$ . Also, the translational speed  $v_z$  is limited by the size of  $\alpha$ .

*rpm* is a 4 dimensional action space. The actions are processed within a range of  $\alpha$  about the *hover\_rpm* to a 4-tupel of rpms. The drone is not limited in any dimensionality, but the task increases in complexity. Due to Equation (2.1), Equation (2.2), Equation (2.3), Equation (2.4) even a small difference in rpms can lead to an unstable flight or even a crash, because the roll or pitch angle is too high. Also, all translational and rotational speeds are limited by the size of  $\alpha$ . Because of the higher complexity, it is expected, that the training takes noticeable more time.

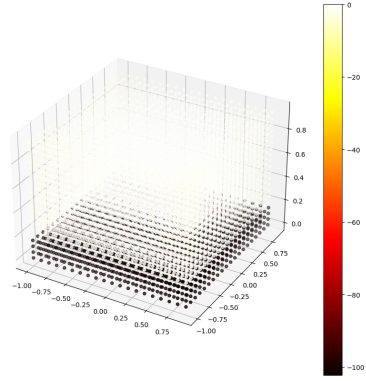
*vel* is a top level, 4-dimensional action space. The action consists of a velocity vector  $v = (a_0, a_1, a_2)$  and its size  $a_3$ . Since it is a top level action space, the actions are not corresponding directly to the rpms, but a pid controller (Section 2.1.2) is used in order to control the rpms. The basic pid controller is part of Gym Pybullet Drones [PZZ<sup>+</sup>21] and must be tuned for the used quadcopter. It mainly receives the state  $S$  of the drone, as well as the targeted velocity, which is calculated with the use of the actions. Therefore,  $a_3$  is multiplied with speed limit of the drone in order to derange the action. Also, the velocity vector  $v$  is normalized to a length of 1.

The drone is not limited in any dimensionality and the task is less complex than setting rpm directly. Stability of the flight is now mainly controlled by the pid controller, so it is bounded by the typical pid constraints in harsh environments and not adaptable.





(a)



(b)

**Fig. 4.2:** Visualization of the used reward function with the goal  $(0, 0, 0.5)$  and a colour scale: the distance dependent part (a) and the big penatalization, because the drone hit the ground (b)

## Reward

## Constraints

#### **4.1.2 WindSingleAgentPathfollowingAviary Environment Class**

## 4.2 Scripts & Evaluation Tools

### 4.2.1 Learning Script

---

**Algorithm 2:** Learning Script

---

- 1 Parse the arguments to the Script
  - 2 Check the parsed arguments on contradiction and raise `ParsingError` if needed
  - 3 Create training environment with the parsed parameters
  - 4 Define the size of the actor and critic policy network
  - 5 Create or load the model based on the parsed load argument
  - 6 Create evaluation environment with the parsed parameters
  - 7 Define evaluation callbacks
  - 8 Learn the model with the parsed amount of time steps
  - 9 Save the model with the parsed preferred name
-

### 4.2.2 Evaluation Script

---

**Algorithm 3:** Evaluation Script

---

```
1 Parse the arguments to the Script
2 Check the parsed arguments on contradiction and raise ParsingError if needed
3 Create evaluation environment with the parsed parameters
4 Define the size of the actor and critic policy network
5 Load the model
6 Instantiate an EvalWriter and evaluate the model
7 if gui parameter was parsed as true then
8     Instantiate a test environment, logger and a pathplotter
9     Reset environment and save an observation  $\sigma$ 
10    Save current time  $t$ 
11    repeat
12        Get action  $a$  and the observation  $\sigma$ 
13        Execute a step in the test environment and receive an observation  $\sigma$ , a
            reward  $r$ , a done value and a info
14        Add current pose to pathplotter
15        Log current drone state
16        Sync simulation time  $t_{sim}$  to the real time  $t$ 
17    until done or time is over
18    Close test environment, show the plotted path and the values.
```

---

### **4.2.3 Evaluation Tools**

#### **EvalWriter Class**

## PathPlotter Class

## **5 Evaluation**

### **5.1 Drone Model**



## 5.2 Setup

## 5.3 Metric

## 5.4 Results

## **6 Conclusion & Future Work**

## 7 Appendix

```
#imports
...

if __name__ == "__main__":
    parser = argparse.ArgumentParser(description="Script that allows to
                                         train your RL Model")

    parseparameters(parser)
    ARGS = parser.parse_args()
    check(**vars(ARGS))
    ...

    # Create training environment
    sa_env_kwargs: dict = dict(aggregate_phy_steps=5, obs=ObservationType
                               ('kin'), act=act, mode=mode,
    total_force=total_force, upper_bound=upper_bound, drone_model=drone,
                               gui=gui,
    debug=debug_env, episode_len=episode_len)
    train_env = make_vec_env(WindSingleAgentAviary, env_kwargs=
                               sa_env_kwargs, n_envs=cpu, seed=0)

    onpolicy_kwargs: dict = dict(activation_fn=torch.nn.ReLU, net_arch=[
                               256, 256])

    if load == DEFAULT_LOAD:
        model = PPO(a2cppoMlpPolicy,
                    train_env,
                    policy_kwargs=onpolicy_kwargs,
                    tensorboard_log='results/tb/',
                    verbose=1)
    else:
        model = PPO.load(load, train_env, tensorboard_log='results/tb/')
    # Create evaluation Environment
    eval_env = make_vec_env(WindSingleAgentAviary, env_kwargs=
                               sa_env_kwargs, n_envs=cpu, seed=0)

    # Train the model
    ... # Callbacks...
    model.learn(total_timesteps=steps, callback=eval_callback,
                log_interval=100)
    # Save the model
    model.save(name)
```

Fig. 7.1: shortened excerpt of the python implementation of the learning script

# Bibliography

- [Bis94] Chris M. Bishop: Neural networks and their applications. *Review of Scientific Instruments*, 65(6):1803–1832, June 1994, 10.1063/1.1144830, ISSN 0034-6748. <https://doi.org/10.1063/1.1144830>.
- [CB21] Erwin Coumans and Yunfei Bai: Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2021.
- [DJS20] Dr. Roland Schwaiger Dr. Joachim Steinwendler: *Neuronale Netze programmieren in Python*. Rheinwerk Verlag, 2020.
- [FWXY20] Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang: A theoretical analysis of deep q-learning. In Alexandre M. Bayen, Ali Jad-babaie, George Pappas, Pablo A. Parrilo, Benjamin Recht, Claire Tomlin, and Melanie Zeilinger (editors): *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of Machine Learning Research*, pages 486–489. PMLR, 10–11 Jun 2020. <https://proceedings.mlr.press/v120/yang20a.html>.
- [KMWB19] William Koch, Renato Mancuso, Richard West, and Azer Bestavros: Reinforcement learning for uav attitude control. *ACM Transactions on Cyber-Physical Systems*, 3(2):22, 2019.
- [MP43] Warren S McCulloch and Walter Pitts: A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- [NBMB07] Derek R Nelson, D Blake Barber, Timothy W McLain, and Randal W Beard: Vector field path following for miniature air vehicles. *IEEE Transactions on Robotics*, 23(3):519–529, 2007.
- [PZZ<sup>+</sup>21] Jacopo Panerati, Hehui Zheng, SiQi Zhou, James Xu, Amanda Prorok, and Angela P. Schoellig: Learning to fly—a gym environment with pybullet physics for reinforcement learning of multi-agent quadcopter control. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021.
- [RHG<sup>+</sup>21] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann: Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. <http://jmlr.org/papers/v22/20-1364.html>.

- [SSA17] Sagar Sharma, Simone Sharma, and Anidhya Athaiya: Activation functions in neural networks. *Towards Data Sci*, 6(12):310–316, 2017.
- [SSS13] P.B. Sujit, Srikanth Saripalli, and J.B. Sousa: An evaluation of uav path following algorithms. In *2013 European Control Conference (ECC)*, pages 3332–3337, 2013, 10.23919/ECC.2013.6669680.
- [SWD<sup>+</sup>17] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov: Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [SZWJ22] Zhong Shi, Fanyu Zhao, Xin Wang, and Zhonghe Jin: Deep kalman-based trajectory estimation of moving target from satellite images. In *2022 IEEE 10th Joint International Information Technology and Artificial Intelligence Conference (ITAIC)*, volume 10, pages 71–75. IEEE, 2022.
- [Tur50] A.M. Turing: *The Imitation Game*. 1950.