# Exploring the fundamentals of Reinforcement Learning through Q-Learning methods

Levin Ceglie

KANTONSSCHULE AM BURGGRABEN

**Abstract**

This work discusses and answers two hypotheses. The first being the following: *One of the most basic solution methods of reinforcement learning, namely tabular Q-Learning, is adequate to solve the Cartpole environment, provided by OpenAI-Gym.* And the second: *Using deep reinforcement learning, an agent can learn to play the game of Snake on a human-level performance.* As stated in the hypotheses themselves, they lie within the field of reinforcement learning. Therefore, this work begins by exploring and explaining the theory underlying reinforcement learning. This work will conduct three different experiments to not only answer the stated hypotheses, but also to improve the understanding of the presented algorithm.

We start by understanding what the concept behind reinforcement learning is and set the foundations for this work with a mathematical framework called the Markov decision process (MDP). The MDP is the base for the different solution methods. The first one being tabular Q-learning. To improve our understanding of said solution method, we start by solving a simple gridworld environment, where we discover the immense importance of parameter optimization.

We apply the same tabular Q-learning algorithm to the Cartpole environment provided by OpenAI-Gym. OpenAI-Gym is a programming library, written in python, for developing and comparing reinforcement learning algorithms. With the parameter optimization method explained in the first experiment (i.e., the gridworld environment), we manage to find a set of parameters, with which the agent is able to solve the Cartpole environment in only 164 episodes. That result confirms the first hypothesis.

Before we tackle the second hypothesis, we look into artificial neural networks. We first discuss their role in reinforcement learning, followed by an overview of their structure. Then we dive deeply into the details of how an artificial neural network functions. We also learn the basics of optimizing an artificial neural network to assume a given function. With that knowledge, we incorporate artificial neural networks into our existing tabular Q-learning algorithm and thus end up with a deep Q-learning algorithm, which allows greater complexity.

Finally, we set up the Snake game as to fit in the framework of an MDP. Most interesting is the state we choose to return to the agent. We find that there are two options when thinking about the state representation given to the agent, a top-down view or a 2D vision, and within those two categories, there are again many different options. We continue by optimizing the agent's parameters as best as possible within a reasonable timeframe. With those optimized parameters, we train sixteen different agents using eight different state representations. Besides the significant performance differences between the state representations, we find that most agents were able to outperform a human test group, consisting of fifteen people, by a very significant margin. To be more specific our best agent averaged a score of 27 on a nine by nine grid, which is equivalent to filling 35% of the grid. Whereas, the human test group achieved an average score of 1.1 that is equivalent

to filling 3% of the grid. Furthermore, our results surpassed ones from literature, as their agent filled 8% of the grid on average. These results confirm our second hypothesis.

# Contents

# 1  Introduction

This Matura Paper is concerned with two hypotheses:

- One of the most basic solution methods of reinforcement learning, namely tabular Q-Learning, is adequate to solve the Cartpole environment, provided by OpenAI-Gym [1].

- Using deep reinforcement learning, an agent can learn to play the game of Snake on a human-level performance.

The motivation behind this Matura Paper can be summed up with the following question. How much can a Matura student with hardly any amount of prior knowledge in the field of reinforcement learning achieve in said field within the scope of a Matura Paper? In order to explore this question and confirm or deny my hypotheses, I will first be concerned with the theory behind reinforcement learning. Following that, I will conduct thorough experiments using different Q-Learning methods.

It is important to mention that this work will only discuss the theory on a high level since it would be too much to cover it in detail within a Matura Paper scope. For a deeper dive into the theory, I can only recommend *Reinforcement Learning, An introduction*, by Sutton and Barto [10].

# 2 Theoretical Overview

Reinforcement learning is two things simultaneously. For one, it is a problem, but it also refers to the solution-methods of said problem. The so-called *Reinforcement Learning Problem* can be defined in one short sentence as the following: *Given an environment, how does one act to maximize a numerical reward signal* [10]. The solution-methods obviously cannot be defined in just one sentence. Section 3 and 4 will be concerned with those. However, for now, we shall focus on understanding what reinforcement learning is all about.

## 2.1 Trial and Error

The fundamental process underlying reinforcement learning is *trial and error* [10], which can often be observed in children playing with toys. There is a well known toy, which serves as a good example. The toy I am referring to is a big cube with holes of different shapes, where the goal is to fit given objects through those holes into the cube.



Figure 1: Example toy [7]

Following the process of trial and error, a child with no prior knowledge would either select an object and try it on every hole or select a hole and try every object on it. If we were now to give the child a reward of any kind for each object it fits into the cube, it would quickly learn which object fits through which hole and adapt its behavior, assuming it fancies the reward. This exact idea is what defines reinforcement learning at its core. The following is a more formal definition, of which the terminology will be explained soon.

*Reinforcement learning (RL) is an area of machine learning concerned with how software agents ought to take actions in an environment in order to maximize the notion of cumulative reward.* [17]

## 2.2 Markov decision process

What we discussed in the previous section shall now be formalized into a more mathematically usable form. In order to do that, we first have to understand certain terms, which are mostly self-explanatory. First of all the *agent*, he is the one who takes *actions* in an *environment*, which may or may not be affected by his actions. The *reward signal* is given to the agent for every action taken. It is a scalar measurement for how well the agent did. Now one might ask, based on what does the agent choose his actions? Well, last but not least, we have the agent's *state* in an environment, which can also be thought of as the agent's observation of the environment. [11] [10]

To give a concrete example, let us look at our situation from the previous section. In this case, the agent would be the child, and the environment would be the toy, including the big cube and objects. The state would be whatever the child observes through his sensory input. In this case, his visual input would be most important. One action might be to push a cube through a square hole into the larger cube, for which the child would receive a positive reward, whatever that might be. Contrary to that, if he had tried to fit a cube through a round hole, he would have either received a negative reward (i.e., a punishment) or a neutral one. At this point, his state might be him seeing a cube in his hand as well as two holes on a large cube in front of him, one shaped like a square the other like a circle.

The just discussed terminology is part of the *Markov decision process* [10] [16] or MDP for short, which provides a mathematical framework to our agent-environment model.



Figure 2: agent-environment model described as an MDP [10]

As visualized in Figure 2 an MDP describes the interactions between an agent and an environment in a sequential order. If we were to form a timeline it would look somewhat like the following: $S_0, A_0, S_1, R_1, A_1, S_2, R_2, A_2, ...$ where the subscript represents the time $t$. As you might have noticed there is no $R_0$, which is due to the fact that a reward always follow an action and since there is no $A_{-1}$, there is also no $R_0$. To further comprehend this notation,

it is important to know that $S$, $A(s)$ and $R$ are the sets of all possible states, actions in a given state $s$ and rewards respectively. An element of those sets would be written as $s$, $a$ and $r$ respectively, thus $s \in S$, $a \in A(s)$ and $r \in R$. To omit confusion I want to quickly indicate the difference between $S_t$ and $s$. $S_t$ is a random variable of the state at time $t$ and $s$ is one of its possible events with the probability $Pr\{S_t = s\}$ of occurring. [10] [11] [9]

However, an MDP is more than just a way to describe a sequence of actions and its consequences. It can also serve as a complete model of the environment's dynamics. Thus it can also probabilistically predict future states and rewards given an initial state and action. Here we introduce the first important function. [10]

$$p(s', r | s, a) \doteq Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\} \tag{1}$$

The function $p$ answers the question: What is the probability of the agent ending up in $s'$ (the successor state of $s$) and receiving a reward of $r$ if he is currently in $s$ and takes the action $a$? In other words $p$ describes the *dynamics* of the MDP. The fact that $p$ fully defines the dynamics of an MDP gives rise to the so-called *Markov property* [10]. If a state includes all the information necessary from past states (i.e., the history) that affects the future states, it is said to have the *Markov property*, which in an MDP every state must have. [9]

## 2.3   Policy and Value Functions

As mentioned in the definition for reinforcement learning in Section 2.1 the agent does not only care about its immediate reward. Instead, he tries to maximize the expected cumulative reward, which is also called the *return* and is denoted $G_t$. [10]

However, before we dive deeper, we have to take a step back and look at the difference between an *episodic task* and a *continuous task*. In an episodic task, there is at least one *terminal state*. If the agent reaches it at any time, the episode is terminated, the environment resets, and he starts over. The time step at which he reaches the terminal state is denoted $T$. Whereas in a continuous task, there are no terminal states. Thus, as given by the name, it is continuous. Hence the environment never resets [10]. From here on out, this work will be assuming an episodic task, as all the conducted experiments fall into that category.

Let us continue by defining the simplest form of the return $G_t$, which is just the sum of all following rewards. [10]

$$G_t \doteq R_{t+1} + R_{t+2} + ... + R_T \tag{2}$$

This form of $G_t$ is rarely used in real applications. A more frequently used definition for $G_t$

is the following:

$$G_t \doteq \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \qquad (3)$$

This definition [10] has two main advantages over the former. First, we added the parameter $\gamma \in [0, 1]$ called *discount factor*, which enables us to chose how farsighted the return should be. If we set $\gamma = 1$, that means an immediate reward is worth as much as a reward $n$ time steps into the future. The contrary happens when $\gamma = 0$. Then we only care about the immediate reward. The ability to adjust $\gamma$ can be especially useful when we do not have a perfect model of our environment, which often leads to uncertainty when looking too far into the future [9]. The second advantage is worth mentioning but does not affect this work since it has to do with continuous tasks. The introduction of $\gamma$ enabled the return to be a finite value (as long as $\gamma < 1$), even if the task is continuous, thus never-ending. In our previous definition of the return in (2), a continuous task would have an infinite undiscounted sum as its return (which is rather hard to work with) since there is no terminal state. [10]

Now as we defined the *return* we can get to *policies* and *value functions*. However, before we define them mathematically, we should understand them as a concept. A *policy*, denoted $\pi$, is what an agent uses to chose an action given a state. More accurately, $\pi$ is a distribution over actions given a state [9]. As a function it would be written as $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$ and return the probability of $a$ being chosen given $s$. On the other hand, a value function, as suggested by the name, outputs the *value* of a given state. The value of a state $s$ is defined as its expected *return*. It can formally be defined as the following: [10]

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t | S_t = s] \qquad (4)$$

A value function is always reliant on a policy. In other words, we cannot know the value of a state if we do not know how we will behave. Thus $\pi$ is found in the subscript. From this first definition, we can easily derive another essential and well-known equation as follows: [10]

$$
\begin{aligned}
v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + ... | S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + ...) | S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
v_\pi(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v(S_{t+1}) | S_t = s] \qquad (5)
\end{aligned}
$$

(5) known as the Bellman Equation gives rise to a recursive property of the value function [10] [9]. To better understand this property, it can be useful to look at a so-called backup diagram shown in Figure 3.
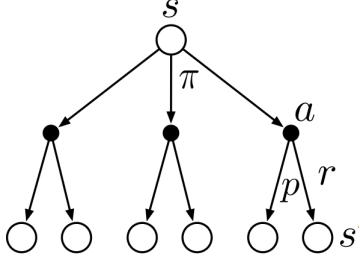
9

Figure 3: Backup diagram for $v_\pi$ [10]

The circles and dots represent states and actions, respectively. Now we can easily see why the recursive relationship is true. The value of $s$ $(v_\pi(s))$ is equal to the immediate reward $r$ plus the value of the successor state $s'$, which is discounted because it is one time step into the future $\gamma v_\pi(s')$. From this diagram it also makes sense, why there is a $\mathbb{E}[\cdot]$ in (5). The reason is that it is not certain what $r$ and $s'$ are going to be. They depend on the agent's policy $\pi$ as well as the environment's dynamics $p$. Therefore we use their random variables $S_{t+1}$ and $R_{t+1}$ in combination with $\mathbb{E}[\cdot]$. [10]

## 2.4 Action-Value Function

The *action-value function* is not too different to the ordinary value function discussed in the last section, which takes the state as an input and returns the value of that state when following a particular policy. On the other hand, the action-value function takes a state and an action as an input and returns the value of that state when first taking the given action and following a specific policy thereafter. [10]

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1}|S_t = s, A_t = a] \tag{6}$$

As you can see (6) is hardly different to $v_\pi$. Another important difference between these two functions is that $q_\pi$ allows us to check which action will result in the biggest estimated return, whereas $v_\pi$ tells us the value if we were to act according to $\pi$. This feature of $q_\pi$ will be quite handy to us later on when it comes to behaving optimally. As one might expect, there is also a Bellman equation for $q_\pi$, which can easily be derived from (6) by removing the expectation. [10]

$$q_\pi(s, a) = \sum_{s',r} p(s', r|s, a)[r + \gamma \sum_a \pi(a'|s')q_\pi(s', a')] \tag{7}$$

Once again, we can form a backup diagram [10] for a better understanding of the recursive property. This time around, it starts with a state-action pair, followed by a state and the
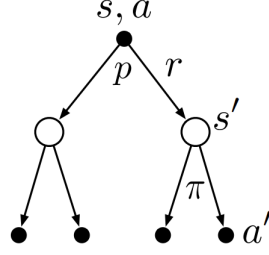
10

different actions available in that state.



Figure 4: Backup diagram for $q_\pi$ [10]

## 2.5 Bellman's Principle of Optimality

Bellman's principle of optimality states [10]: There is always at least one policy that is better or equal to all other policies, called the *optimal policy*. However, what does it mean for a policy to be better than another policy? Well, since our goal is to maximize return, we can compare the corresponding value functions. To demonstrate that, let us assume two policies $\pi$ and $\pi'$, where the latter is known to be better. So if $\pi' > \pi$ then $v_{\pi'}(s) > v_\pi(s)$ must hold for all $s \in S$ [10]. This is saying that the expected return must be greater for all states to call a policy better than another. As stated above, there is at least one optimal policy, which means that there can be multiple optimal policies, which we denote $\pi_*$. Them being equal to each other and better than every other policy leads to a shared *optimal value function*, defined as the following:

$$v_*(s) \doteq \max_\pi v_\pi(s), \qquad \forall s \in S \tag{8}$$

Similarly a shared *optimal action-value function* exists, defined as:

$$q_*(s, a) \doteq \max_\pi q_\pi(s, a), \qquad \forall s \in S \land \forall a \in A(s) \tag{9}$$

[10]. Having defined the optimal action-value function $q_*$ we can now easily define the optimal policy because given a state, we can compare the values of the available actions using $q_*$ and pick the one with the highest value. Mathematically we can write that as the following:

$$\pi_*(s) = \arg\max_a q_*(s, a) \tag{10}$$

In section 2.3 we derived the Bellman equations, which made the recursive property of $v_\pi$ and $q_\pi$ apparent. Similar to them there exist the *Bellman optimality equations*, which give rise to

11

the recursive properties of $v_*$ and $q_*$. [10] [9]

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1})|S_t = s, A_t = a] \tag{11}$$

$$q_*(s,a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a')|S_t = s, A_t = a] \tag{12}$$

Once again, we can form a backup diagram to represent the equations graphically
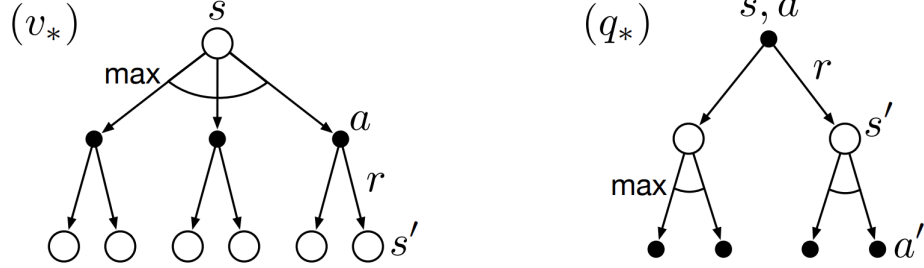


Figure 5: Backup diagram for $v_*$ and $q_*$ [10]

As one can see, these are nearly identical to our previous backup diagrams for $v_\pi$ and $q_\pi$. The only thing that has been added are the arcs, which indicate that the optimal action (i.e., the one with the highest expected return) is being taken, rather than following a given, possibly suboptimal, policy.

# 3 Tabular Q-Learning

This section will be looking at a first solution method, called *Tabular Q-Learning.* To do so, we will first discuss the algorithm itself and then actually use it in two different experiments. By doing so, we will get a better insight into how the algorithm works and what one has to consider when implementing it.

## 3.1 The Algorithm

The Q-Learning Algorithm was introduced by Watkins in 1989 [14], which is to this day still seen as a significant breakthrough in the field of reinforcement learning. The goal is to approximate $q_*$ by incrementally updating the action-value function $Q$. This is done by the following update rule:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)], \tag{13}$$

where $\alpha$ is the *learning rate*, sometimes also referred to as *step-size*, which defines how drastically we update our estimates. The learning rate is limited to an interval, $\alpha \in \, ]0, 1]$. The smaller $\alpha$ is, the more conservative the learning is. It will rely more on old estimates rather than adopting new information. Setting the learning rate to 1 is the complete opposite. It will overwrite its old estimate. If we rearrange (13) just a little bit, it will become clearer why this is true.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)] \qquad \text{(from (13))}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a)] - \alpha Q(S_t, A_t)$$

$$Q(S_t, A_t) \leftarrow (1 - \alpha) \cdot Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a)] \tag{14}$$

As visible in (14) we take $(1 - \alpha)$ of our old estimate and combine it with $\alpha$ of the term in brackets, which at first might not look all too familiar. But if compared to the Bellman optimality equation for $q_*$ (12), we can see where the update rule stems from. So the $Q$-function will continue to be updated using just one value that is known to be right $(R_{t+1})$, since the agent, receives it as feedback from the environment until it has completely converged to $q_*$ because then everything left from the arrow is precisely equal to everything on the right side. Note that updating, not only using actual experience but also using current estimates, is called bootstrapping [10]. One value of $Q$ is updated using another value of it. Furthermore, this update rule does not rely on the current policy being followed by the agent. This property is referred to as off-policy learning [10] and it can be quite useful when implementing the

13

algorithm because it does not matter how the agent behaves. As long as every state-action pair is visited often enough, $Q$ will converge to $q_*$. Watkins proved its convergence in 1992 under certain conditions [13]. However, due to time limitations, this work will not cover the how and why of the convergence.

With this update rule and enough data (experience), the agent can approximate $q_*$, which he can then use as a policy by always picking the action with the highest value (10). In other words, if one knows $q_*$ he knows $\pi_*$. But in order for an agent to obtain enough data, he must also have a policy to operate in a given environment. One could argue to always act greedy with respect to $Q$, in other words, to always pick the action with the currently highest estimated value. However, when looking at a simple example, it becomes clear why that might not be the best choice.

Imagine you are visiting a city for the first time, and you have to choose a restaurant for every meal since you cannot cook. Your initial estimates of the restaurants are neutral $(Q(\cdot) = 0)$, so you randomly chose one of them for the first meal. It happens to be a pleasant experience, and you update your estimate for choosing that restaurant positively. Following the greedy policy with respect to your current estimates $(Q)$, you would continue to pick that same restaurant indefinitely. However, for you, there is no knowing whether there is a restaurant of higher value.

The problem demonstrated by the example is that if an agent was to always *exploit* its current estimates and never *explore* other actions, it would often get stuck on a suboptimal policy and action-value function. A well-known method to address this problem is to use an $\varepsilon$-greedy policy during training [10]. When picking an action using $\varepsilon$-greedy, the action will be chosen at random with probability $\varepsilon$ and with probability $1-\varepsilon$, the action will be greedy with respect to the current $Q$. Thus there will always be some amount of exploration depending on the parameter $\varepsilon \in \, ]0, 1]$.

When using the *Tabular* Q-learning algorithm, the $Q$-function will be represented by a look-up table (i.e., array) of size $S \times A$. Thus it will contain a numerical value (current estimate) for each possible state-action pair and is usually referred to as the $Q$-table [10].

| | $a_0$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | $a_6$ | ... | $a_n$ |
|---|---|---|---|---|---|---|---|---|---|
| $s_0$ | | | | | | | | | |
| $s_1$ | | | | | | | | | |
| ... | | | | | | | | | |
| $s_m$ | | | | | | | | | |

Figure 6: $Q$-table for the Tabular Q-learning algorithm

The number of actions and states in Figure 6 have no meaning at all. They were chosen arbitrarily to fit nicely. The Figure's sole purpose is to give a graphical representation of the $Q$-table. In every empty cell belongs the current estimate of the according state-action pair. So every time the $Q$-function is called, the value is just looked up in the $Q$-table. Having discussed all that, here is the complete algorithm that was used in the following experiments.

---
**Algorithm 1** Tabular Q-learning algorithm [10]
---
Parameters: learning rate $\alpha \in \, ]0, 1]$, exploration rate $\varepsilon \in \, ]0, 1]$

Initialize $Q(s, a)$, for all $s \in S, a \in A(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

**for** each episode **do**

    Initialize S

    **for** each step of episode **do**

        Choose $A$ from $S$ using policy derived from Q using $\varepsilon$-greedy

        Take action $A$, observe $R$, $S'$

        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

        $S \leftarrow S'$

    **end for** if $S$ is terminal

**end for**

---

## 3.2 Gridworld

The primary purpose of this experiment is to gain a better understanding of how the Q-learning algorithm works. Also, it will become clear how vital the adjustment of the different parameters can be.

### 3.2.1 The Environment

I created a relatively simple gridworld environment that looks like the following:
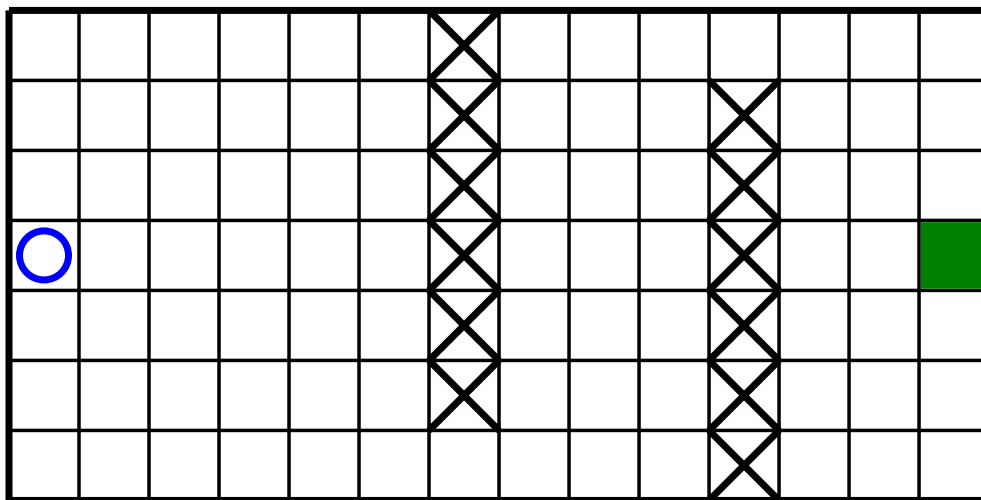


Figure 7: Gridworld environment

In Figure 7 we see the initial state of the environment, where the blue circle is the agent, whose goal is to reach the green rectangle while avoiding the walls represented by the crossed cells. He can move up, down, right, and left, with a reward of -1 on every transition except on the following:

- trying to step out of the grid will result in no movement and a reward of -10
- stepping into the walls will result in death (i.e. termination of the episode) and a reward of -100
- stepping into the green rectangle will also terminate the episode, but with a regular reward of -1 for the transition.

In this setting, the shortest paths between the agent's starting cell and the green cell will result in a total undiscounted reward of -25. Note that there is not one shortest path but multiple. For example, in order to pass the first wall, it does not matter whether you first adjust your vertical position or your horizontal one. There is a shortest path for both options. Contrary to that, there is only one optimal (i.e., highest) score you can achieve, which, as

mentioned earlier, is -25. The beauty of this example is that Bellman's principle of optimality (Section 2.5) suddenly becomes quite apparent. There are multiple optimal policies (paths), but there is just one optimal value and action-value function since all optimal policies end up achieving the same optimal score. To clarify, by score, I mean the sum of undiscounted rewards achieved during an episode.

The next important thing is the state-space of the environment. This can significantly impact performance, especially when using a tabular solution method, as the next experiment will demonstrate. Luckily in this experiment, our gridworld is already limited to $7 \cdot 14 = 98$ states because only that many cells are available to the agent. There are multiple ways of how the agent might observe these different states, but the simplest one is probably as coordinates. Since we will be using a 2-dimensional array (i.e., matrix) to represent the environment, it is most comfortable for the coordinates to have the form (row, column) starting with cell $(0, 0)$ in the upper left corner. The agent's goal would thus be located at $(3, 13)$.

### 3.2.2 Implementation of the Algorithm

Having defined our environment's most important properties, we can get to solving it by implementing Algorithm 1. Since this is a rather simple environment, our goal should be for the agent's $Q$-table to converge to $q_*$ and do so in as few episodes as possible. The only difference one can make as the developer is adjusting the different parameters of the algorithm and finding the best ones. Before we start our search, it is important to note that there are actually two additional parameters than the algorithm initially suggested. The discount rate defines how farsighted the agent should be, but due to the simplicity and lack of uncertainty in this environment, we will always set $\gamma = 1$. And then, there is the initial value of each state-action pair, which can be chosen arbitrarily except for the terminal states. From this point forward, I will refer to that value as $q0$. To clarify, by performance, at least in this simple experiment, I mean the number of episodes it takes the agent's $Q$-table to converge to $q_*$. The fewer episodes it took, the higher the performance. Using the $\varepsilon$-greedy policy, it is not wise to track the score achieved by the agent on each episode and call that his performance. That is due to the fact that with probability $\varepsilon$ he will choose a random action, which can lead to a suboptimal score, even though his $Q$-table might already have converged to $q_*$. But in this simple environment we can easily calculate the optimal state-action values for the starting state ($s_0 = (3, 0)$). So we could compare them to the agent's current estimates for $s_0$. After some testing, I found that if they are equal ($Q(s_0) = q_*(s_0)$), it is pretty safe to say that the $Q$-table has converged close enough to $q_*$. Note that the $Q$-table's values will only be calculated to eight decimal places and then rounded. To demonstrate that, I let an agent

17

with arbitrary parameters play until $Q(s_0) = q_*(s_0)$ and at that point, I saved his $Q$-table and rendered his greedy policy with respect to his $Q$. Here are the results:
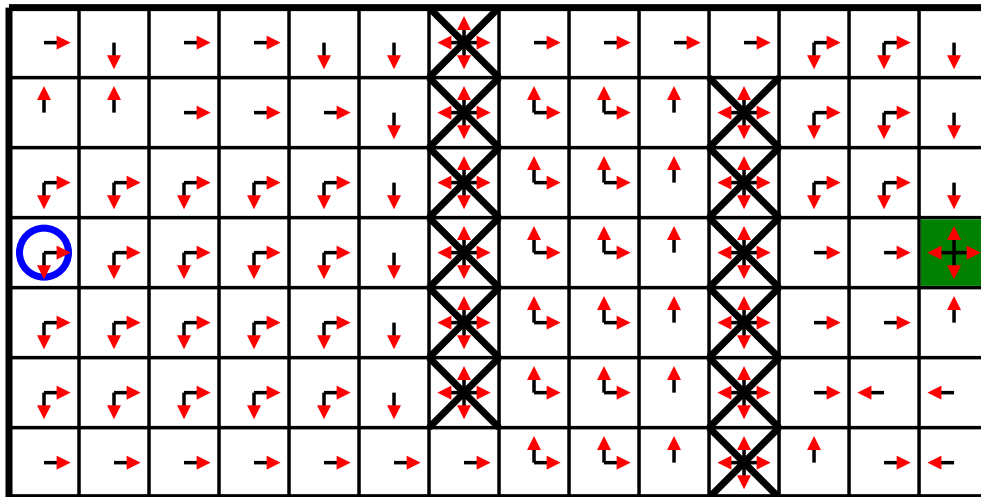


Figure 8: Gridworld, greedy policy after 4490 episodes, $\varepsilon = 0.1, \alpha = 0.5, q0 = 0, \gamma = 1$

The red arrows represent the action with the highest estimated value. If there are two in one cell, then they share the same value. With the exception for the upper left and lower right corners, the agent's $Q$-table did converge to $q_*$ and thus he has found $\pi_*$. And since these corners are never visited, when following $\pi_*$ they are negligible. To conclude we can assume that the agent's $Q$-table has converged close enough, as soon as $Q(s_0) = q_*(s_0)$.

### 3.2.3   Parameter Optimization

As mentioned earlier, our goal is to find the best set of parameters. To do so, we start with arbitrary ones and then adjust one parameter at a time to see how it affects performance. Let us first take a look at $\varepsilon$. After choosing different values for $\varepsilon$, I ran 20'000 episodes 100 times and averaged the performance over those 100 times for each value of $\varepsilon$.
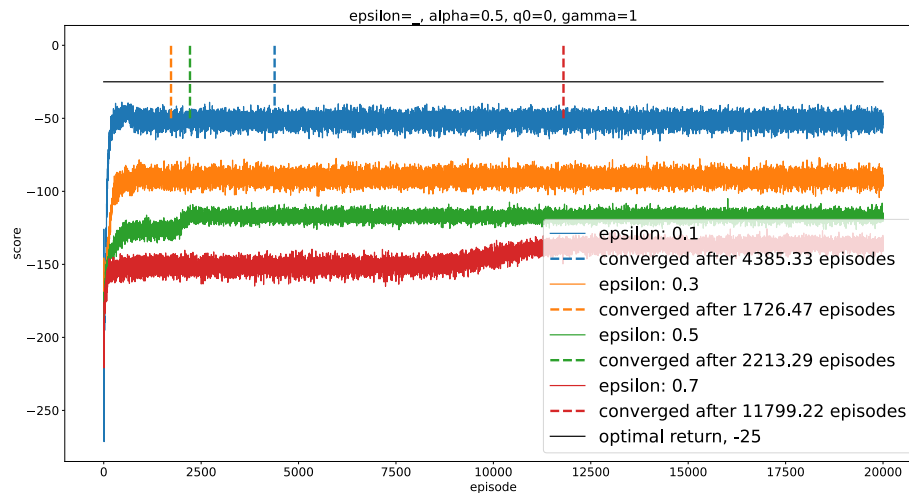
Figure 9: Gridworld, effects of adjusting $\varepsilon$

As visible in Figure 9, there are always 2 different metrics plotted for each $\varepsilon$ (color). The dashed vertical lines mark the episode at which the $Q$ converged close enough to $q_*$ as discussed earlier, and other lines show the score of each episode during training. Note that the score achieved during training is not crucial. More important is the amount of episodes it took for the $Q$-table to converge, which is always clearly visible in the legend. Moreover, I also plotted a horizontal black line for reference, which shows the optimal score of -25. The results suggest a sweet spot around $\varepsilon = 0.3$ at which the $Q$ converges fastest. Note the difference in scores, which is a direct result of a higher value for $\varepsilon$, since the agent is only allowed to act optimally with probability $(1 - \varepsilon)$.

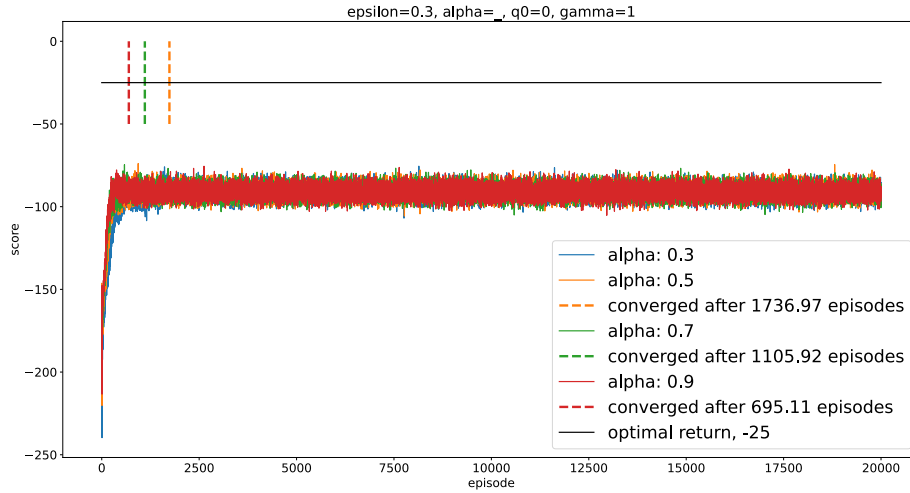Next up I adjusted $\alpha$ and once again ran 20'000 episodes 100 times for each value of $\alpha$, here the results:

Figure 10: Gridworld, effects of adjusting $\alpha$

The first thing to notice is that with $\alpha$ set to 0.3, the $Q$-table did not manage to converge close enough within 20'000 episodes, but as we increase $\alpha$, the performance also increases, which suggests that if we had set $\alpha$ even higher than 0.9, the performance would have increased as well. However, for this experiment, we do not necessarily seek the absolute best performance. Therefore we will stick to $\alpha = 0.9$.

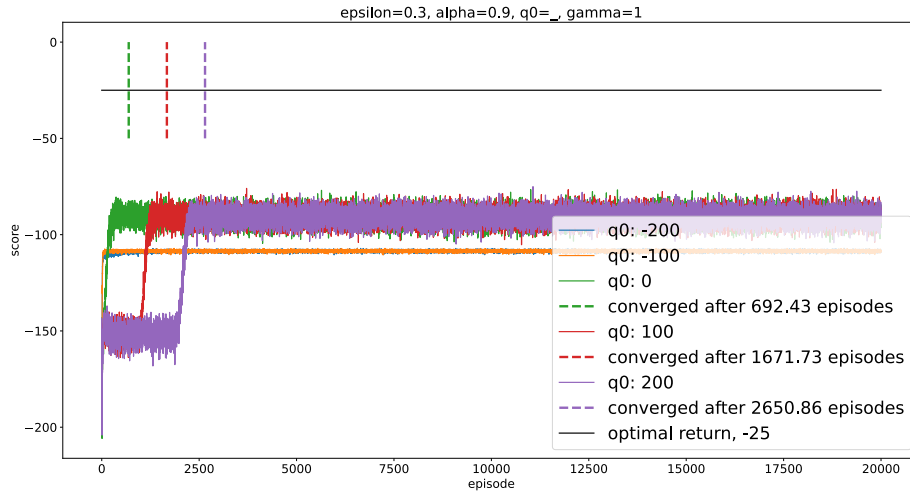Last but not least let us adjust $q0$. Again I ran 20'000 episodes 100 times for each value of $q0$.



Figure 11: Gridworld, effects of adjusting $q0$, I

After I saw this plot, I was unsure why $q0 = -100$ and $q0 = -200$ did not converge at all. My first thought was, if the initial estimates were too close to a worst-case scenario, e.g.,

walking right until he steps into the wall, the agent would see the worst-case scenario as a good solution. At that point, it is up to the 0.3 ($\varepsilon$) probability of taking a random action to completely avoid the walls and reach the green rectangle often enough, to realize that it is a better solution. Therefore I reran it with different values for $q0$.



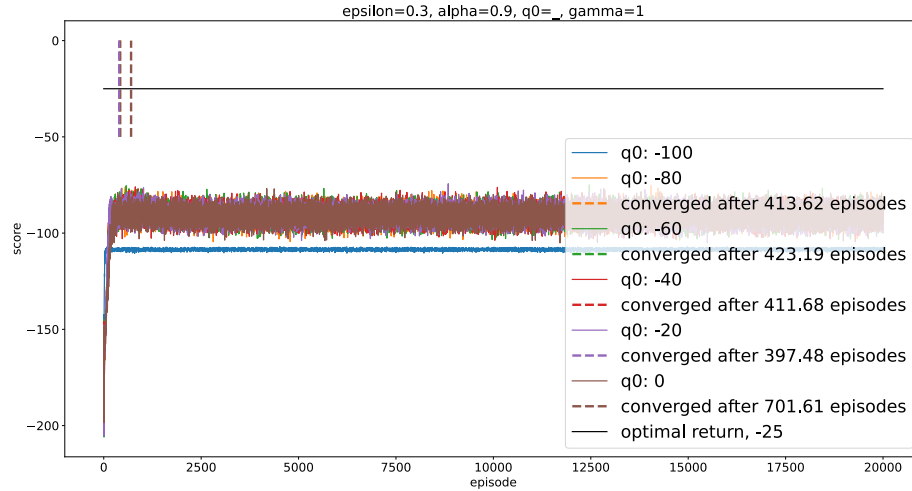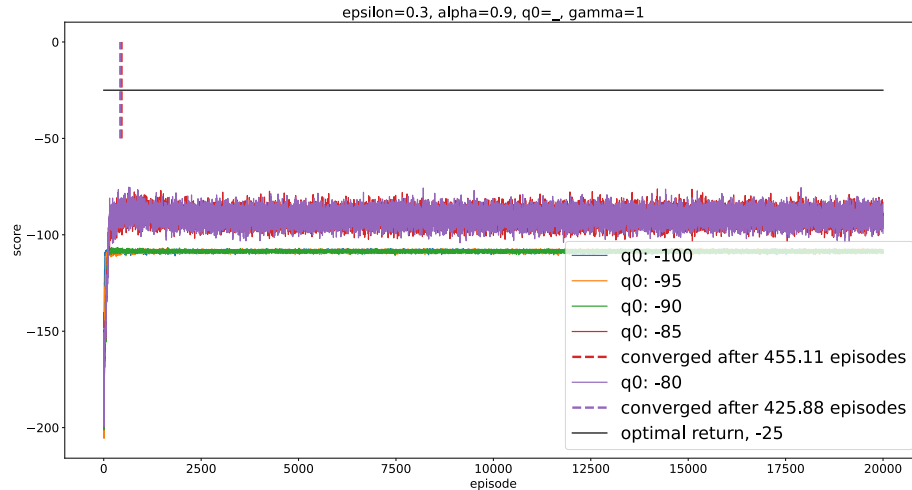Figure 12: Gridworld, effects of adjusting $q0$, II



Figure 13: Gridworld, effects of adjusting $q0$, III

After analyzing Figure 12 and 13, my earlier stated hypothesis might be true. To investigate this further, I once again rendered the greedy policies for $q0 = -85$, were it on average converged after about 455 episodes, and also for $q0 = -90$, where it did not converge within 20'000 episodes.

21

Figure 14: Gridworld, greedy policies, top: $q0 = -85$ after 430 episodes, bottom: $q0 = -90$ after 20'000 episodes, both: $\varepsilon = 0.3, \alpha = 0.9, \gamma = 1$

As visible in Figure 14, for $q0 = -85$, the greedy policy with respect to $Q$ after only 430 episodes is an optimal one, whereas for $q0 = -90$, it is clearly suboptimal and it prefers to walk straight into the walls. This loosely confirms my hypothesis. To sum up, we can safely say that having initial estimates below or near a worst-case scenario can lead to a suboptimal policy and should be omitted in this environment. Thus, to be on the safe side, one should start with somewhat optimistic estimates.

### 3.2.4   Takeaway

To conclude this experiment, I want to emphasize how crucial it is to optimize the parameters to the best of one's ability. It really can make a significant difference in the learning

performance, as we saw with $q0$ in this specific environment.

## 3.3 Cartpole

This section will be concerned with my first hypothesis, which states: *One of the most basic solution method of reinforcement learning, namely tabular Q-Learning, is adequate to solve the Cartpole environment, provided by OpenAI-Gym [1].*

The Cartpole environment consists of a pole, which is attached to a cart. The goal is to prevent the pole from falling over. To do so, there are two actions available to the agent, either pushing the cart to the right or to the left.



Figure 15: Cartpole environment provided by OpenAI-Gym [1]

The state-space of this environment was described by OpenAI [1] as follows:

Table 1: Cartpole State-Space

| Num | Observation | Min | Max |
| --- | --- | --- | --- |
| 0 | Cart Position | -4.8 | 4.8 |
| 1 | Cart Velocity | -Inf | Inf |
| 2 | Pole Angle | -0.418 rad | 0.418 rad |
| 3 | Pole Angular Velocity | -Inf | Inf |

As we can see it consists of four numbers, which each have a specific range. Hence we are dealing with a continuous state-space, which cannot be represented in a finite $Q$-table. Thus

we somehow have to transform the given continuous state-space into a discrete one. This is actually less complicated than it might sound. Let us take the Cart Velocity as an example, which *theoretically* has a range from -inf to inf. But after letting an agent, who picks actions at random play 2500 episodes, I created the following plot containing the observed Cart Velocity at each time step.



Figure 16: Cartpole, observed Cart Velocity

As one can see, in practice, the Cart Velocity only really ranges from -2 to 2. In addition to that, we have to discretize the received observation. We do that by splitting up the range from -2 to 2 into smaller chunks. For example, we say everything between -2 and -1.8 is put into the first chunk (0), everyth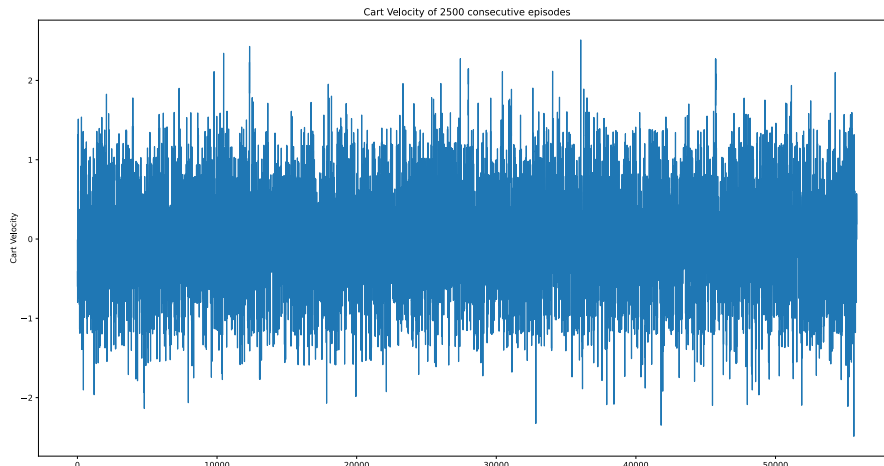ing between -1.8 and -1.6 into the second chunk (1), and so on. By doing that, we end up with a finite state space, which can then be represented in an agent's $Q$-table. Note that the number of chunks we divide the continuous space into is also a parameter, as well as the smaller adjusted range for each observation. Also note that when implementing it, the first chunk would actually range from -inf to -1.8 to omit any errors, as the Cart Velocity could theoretically go below -2. The same goes for the last chunk, which would range from 1.8 to inf.

However, before we dive any deeper into actually solving the environment, we have to define what it means. OpenAI considers it solved if the average return is greater than or equal to 195.0 over 100 consecutive trials [1]. Their requirements do not make it clear whether it has to be achieved during training or in testing, because in training the agent acts $\varepsilon$-greedy, whereas in testing we let him chose his actions greedily. I decided that the agent must achieve it during training.

In the gridworld experiment, we used a fixed value for $\varepsilon$ during training. Thus the agent

will always explore at the same rate during training. However, a well known and also by tests of mine shown to be effective strategy is to let the agent explore a lot initially, and overtime as he gets better, we slowly lower the exploration rate ($\varepsilon$) [4]. This is actually a relatively intuitive strategy when you do not know anything about your environment, it makes sense just to try a bunch of things, and as you get to know it better, you can start to actually do things that you expect to result in a high reward. Another welcomed advantage is that by doing so, we will be able to achieve our required average return during training because as we saw in Figure 9 the achieved score during training is highly affected by $\varepsilon$.

After following the same procedure as in the gridworld experiment to optimize the parameters, I came up with a set of parameters, with which the agent was able to solve the environment in only 164 episodes.



Figure 17: Cartpole, solved after 164 episodes

As visible in the plot above the agent did successfully solve the environment. The Following parameters were used (in Source Code E) to achieve the results.

- obs_high=[2.4, 2, 0.21, 1.8]
- obs_low=[-2.4, -2, -0.21, -1.8]
- obs_chunks=[1, 1, 12, 12]
- q0=0
- gamma=1
- epsilon_function=1/(episode+1)
- alpha=0.1
- S_T_reward=-2000

There are two more things I need to mention. First, there is a parameter called S_T_reward,

which I added to give the agent additional punishment when he fails. So every transition into $S_T$, the terminal state, is punished with a reward of -2000. The second thing worth mentioning is the function I used to slowly decrease epsilon over time. When plotted, it looks like the following.



Figure 18: Cartpole, decreasing epsilon

The agent starts by only exploring, but he then very quickly decreases his exploration rate. After 100 episodes, his probability to explore is below 0.01, which is vital to solving the environment since we want him to reach the threshold during training. If the exploration rate was, for example, 0.1, he might fail every $10^{th}$ episode, which would lead to not reaching the required threshold.

With these results, I can confirm my first hypothesis, which states: *One of the most basic solution methods of reinforcement learning, namely tabular Q-Learning, is adequate to solve the Cartpole environment, provided by OpenAI-Gym [1].*

# 4 Deep Q-Learning

In the previous section, we used a simple look-up-table to represent the $Q$-function. This approach worked fine for the former simple experiments since they had a rather small state-action space. For the Gridworld the look-up-table had a size of $7 \times 14 \times 4 = 392$ and for the Cartpole experiment it was $12 \times 12 \times 2 = 288$. In this section, the goal will be for an agent to learn how to play Snake, which will have a way bigger state-action space. Therefore we are switching from a tabular to a function approximation method, namely Deep Q-Learning. This will allow for greater complexity and also reduce the computation needed. [8] [9, Lecture 6]

## 4.1 Neural Networks

The currently most popular function approximator for reinforcement learning must be the neural network. Note that by neural network, this work refers to the artificial one and not to the biological one.

### 4.1.1 The Structure of a Neural Network

A neural network is essentially just a function that can be scaled to arbitrarily many input and output dimensions. It is often thought of as very complicated, but as we start to break it down, one will quickly realize that it is easy to understand how the function works. Let us first look at a small neural network to understand what it does. Note that this work only considers linear layers for the neural network.

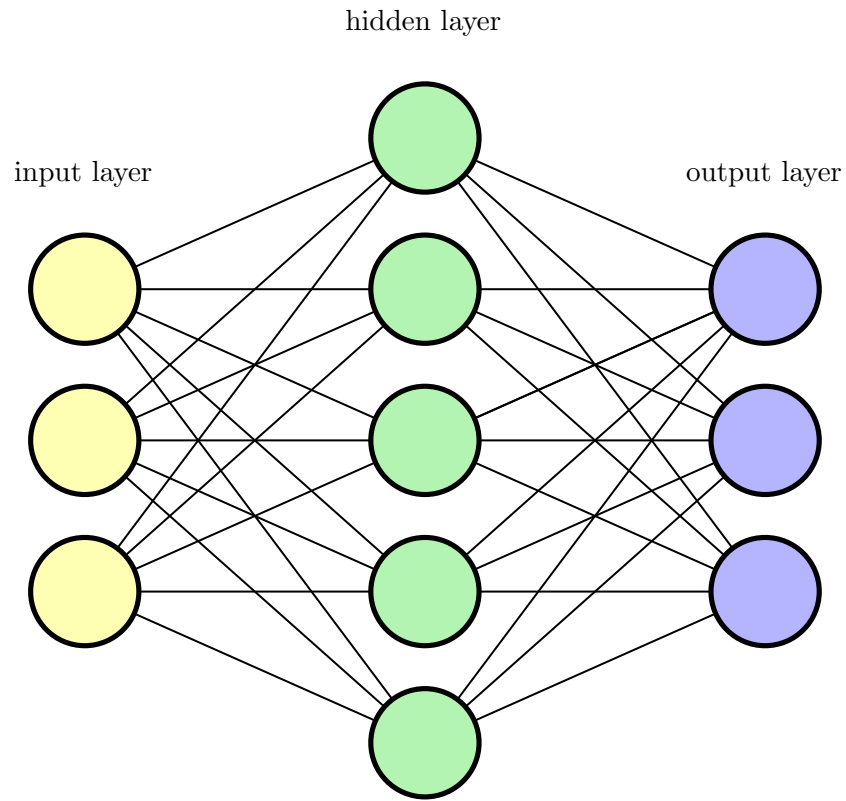Figure 19: Simple neural network

Figure 19 shows a simple neural network consisting of three layers: an input, a hidden, and an output layer. To understand how a forward pass through the neural network works, we will zoom in on one neuron of the hidden layer.
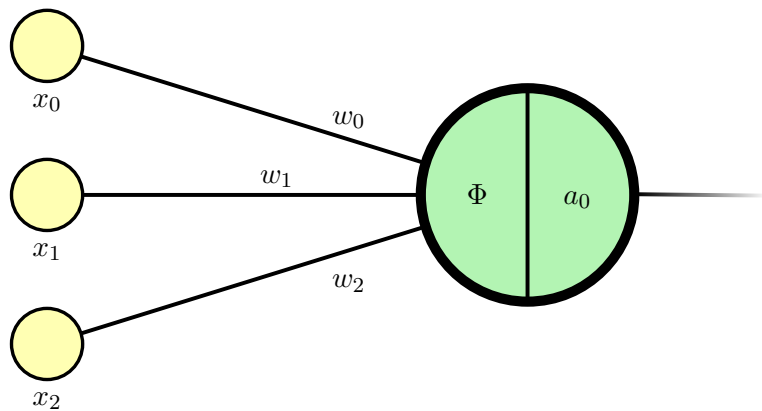


Figure 20: Zooming in on one neuron

Visible in Figure 20 we can see said neuron, which has three incoming values $x_0$, $x_1$ and $x_2$,

which all have an according weight $w_0$, $w_1$, $w_2$. Now the first thing that is calculated as these values are passed to the neuron is $\Phi$. That is done by multiplying every input value by its according to weight and adding a bias $\beta$ afterward. To write it down formally, it would look like the following:

$$\Phi(\vec{x}, \vec{w}) = \sum_i x_i w_i + \beta = \langle \vec{x}, \vec{w} \rangle + \beta \tag{15}$$

As one can see, there are two ways to write $\Phi$. The first is just a simple summation as described above, but in the second one, we build the dot product of the vectors $\vec{x}$ and $\vec{w}$ and add a bias $\beta$ afterward, which will make things easier. In this example, these two vectors would be defined like this:

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}, \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

Here I want to define one more, later important, vector $\theta$ as follows:

$$\vec{\theta} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \beta \end{bmatrix}$$

This vector contains all the parameters of the neural network. In this case, there are only four, as we are only working with a specific neuron of the neural network and only considering what happens there, but the concept and the math behind it could quickly be scaled up to the whole network in Figure 19.

### 4.1.2 Activation Function

Now we know how to calculate $\Phi$ for a certain neuron, but there is one more step to get to $a_0$, the actual value that the neuron passes on. We have to use an *activation function* on $\Phi$. The primary reason for that is to add non-linearity to the neural network. Without a non-linear activation function, all of the linear layers, no matter how many there are, of a neural network would collapse into a single linear layer. Furthermore, non-linearity can be extremely important when approximating a more complicated function. Some functions simply cannot be approximated well through a linear approximator. Note that we will not be using any activation function on our output neurons, as they are the $Q$-values for a given state. And if we were to use an activation function on them, we would limit their ability to accurately approximate the $Q$-value as an activation function often squishes or cuts off a given input.

There are many activation functions to choose from. The rectified linear unit or ReLU for short is a prevalent choice and was also proven to work rather well [2]. Moreover, it is a straightforward function. It merely takes the maximum between 0 and the given input. If one were to plot it, it would look like the following.



Figure 21: The ReLU activation function

Formally we can define it like this,

$$R(x) = \begin{cases} x & \text{if } x > 0, \\ 0 & \text{otherwise.} \end{cases}$$

or even simpler,

$$R(x) = \max(0, x) \tag{16}$$

We can now finish our calculation of $a_0$ by simply passing $\Phi$ through $R$, thus:

$$a_0 = R(\Phi(\vec{x}, \vec{w})) = \max(0, \langle \vec{x}, \vec{w} \rangle + \beta) \tag{17}$$

### 4.1.3 Optimization Algorithm

Now we know how a neuron transforms input signals into a single output signal. With that, we essentially know how a forward pass through a neural network works, because the same calculations happen in every neuron. But as mentioned earlier, our goal is to approximate a given function and to do so, we rather obviously need some way to adjust the parameters $\vec{\theta}$ so that the accuracy of the neural network, when compared to the given function, improves. Note that said function is usually unknown. Only data points are given. To improve, we first have to know how much, if at all, the approximation is off. We need a so-called *loss*

*function* that gives us a numerical value to measure the accuracy. Once again, there are many to choose from. We will simply take the squared error as our loss function.

For the following we have to define one more term, $\hat{y}(\vec{x}, \vec{\theta})$ the output vector of a neural network dependent on the input $\vec{x}$ and the parameters $\vec{\theta}$.

With that we can now properly define our loss function as follows.

$$L(\vec{\theta}, \vec{x}, \vec{y}) = \sum_i (\hat{y}_i(\vec{x}, \vec{\theta}) - y_i)^2, \tag{18}$$

where $\vec{y}$ is the target output vector for our input $\vec{x}$. By squaring the difference between the approximated output and the target output, $L$ will always return a positive number and return 0 if the neural network perfectly approximated the target function for the given input. Hence we want to choose the parameters $\vec{\theta}$ to minimize $L$. If we were to form the gradient $\nabla_{\vec{\theta}} L(\vec{\theta}, \vec{x}, \vec{y})$, which is the vector of steepest ascent w.r.t. $\vec{\theta}$ or in other words how do we have to adjust our parameters $\vec{\theta}$ as to most quickly increase $L$. A logical thought would be to simply take the negative value of the gradient, $-\nabla_{\vec{\theta}} L(\vec{\theta}, \vec{x}, \vec{y})$, and adjust our parameters with that. This method is known as stochastic gradient descent (SGD). For any given data point $(\vec{x}, \vec{y})$ it can perform a parameter update to lower the loss. Formally the update rule for $\vec{\theta}$ would look like this:

$$\vec{\theta} \leftarrow \vec{\theta} - \eta \cdot \nabla_{\vec{\theta}} L(\vec{\theta}, \vec{x}, \vec{y}) \tag{19}$$

Here we reintroduce a learning rate under a different variable $\eta$ since its implications differ slightly from the one in Section 3. Here it regulates by which amount the parameters follow the steepest descent from their current position. The adjustment of the learning rate can have a massive impact on performance. That is because all we are trying to achieve with the SGD is to find a good enough local minimum in our loss function since it is, in practice, not feasible to land in the global minimum. However, in order to end up in a decent local minimum, our learning rate cannot be too high, as we would end up overshooting all the local minima, but it also cannot be too low for two reasons: For one, the lower the learning rate, the longer it will take to reach a local minimum and second it might get stuck in a bad local minimum.

Figure 22: implications of adjusting $\eta$

Figure 22 shows how adjusting $\eta$ affects the search for a local minimum. In blue, we see a terrible choice for $\eta$, it is way too high, and it is just jumping around like crazy. In black, on the other hand, we see a close to optimal choice for $\eta$. This is obviously all with respect to the given starting point for the parameter. In this case, we only have a single one starting at 5.6. Note that the parameters of the neural network will be initiated mostly at random, and also rather obviously the dimensionality is not to compare. The quite small neural network in Figure 19 already has 38 different parameters to adjust. Thus, its loss function's input space would be 38-dimensional, which is beyond us humans to imagine.

## 4.2   The Algorithm

One of the first to combine reinforcement learning with neural networks was Gerald Tesauro in 1995 with TD-Gammon [12]. However, probably more significant was the breakthrough from DeepMind when they managed to play several Atari games on a superhuman level through deep reinforcement learning, the combination of deep neural networks and reinforcement learning algorithms [8]. The algorithm I used to play snake only differs slightly from Algorithm 1 and is also very similar to the one suggested by DeepMind [8]. The apparent difference is that we no longer use a look-up-table but rather a neural network to represent our $Q$-function. Furthermore, in the previous Tabular-Q-Learning algorithm, we used every transition played by the agent only once to update the $Q$-table. This is very data inefficient and problematic due to the heavy correlation between transitions that happen one after another. Therefore this algorithm uses an experience replay mechanism [8] [6] which stores every transition ex-

perienced by the agent in $\mathcal{D}$ for later use. On every step, a specific amount $(n)$ of data points are randomly sampled from the thus far collected data (agent's memory) $\mathcal{D}$ and are then individually used to optimize the parameters of the neural network by using *Adam* [3], an algorithm for stochastic optimization. This algorithm is essentially a better, more optimized version of the previously introduced SDG method. If $n$ data points are sampled, there will be $n$ gradient descent parameter updates per step.

One more thing to mention is the design of the neural network. The size and amount of the hidden layers can be chosen freely, but the dimensionality of the input layer has to match the observed state's dimensionality. Furthermore the output layer must have a neuron for every possible action, or in other words, the output layer will consist of the $Q$-values for the individual action. This architecture for a neural network is referred to as a Q-Network [8].



Figure 23: Q-Network architecture, observed state as input and $Q$-values as output (size and amount of hidden layers may vary)

Having discussed all of that, I can now present the Deep Q-Learning algorithm, which I used in the following experiment.

**Algorithm 2** Deep Q-Learning algorithm [10] [8] [3]

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize $Q$-function with random parameters $\vec{\theta}$

**for** each episode **do**

    Initialize episode, $s_0$

    **for** t=0, T **do**

        With probability $\varepsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q(s_t, a; \vec{\theta})$

        Take action $a_t$, observe $r_{t+1}$, $s_{t+1}$

        Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in $\mathcal{D}$

        Sample $n$ random transitions from $\mathcal{D}$

        **for** each transition $(s_i, a_i, r_{i+1}, s_{i+1})$ **do**

$$\text{Set } y_i = \begin{cases} r_{i+1} & \text{for terminal } s_{i+1} \\ r_{i+1} + \gamma \max_{a'} Q(s_{i+1}, a'; \vec{\theta}) & \text{for non-terminal } s_{i+1} \end{cases}$$

            Perform gradient descent step on $(y_i - Q(s_t, a_t; \vec{\theta}))^2$

        **end for**

        $s_t \leftarrow s_{t+1}$

    **end for**

**end for**

Notice that for the loss, we only consider the $Q$-value of the action that was actually taken in that transition and not the whole output vector. That is due to the fact that we only know the immediate reward that followed said action. Thus we can only know the target value for that specific action.

## 4.3 Snake

This section will address my second hypothesis, which states: *Using deep reinforcement learning, an agent can learn to play the game of Snake on a human-level performance.* In order to test this hypothesis, I implemented Algorithm 2 on a self-made game environment. Therefore this section will first elude to the most important aspects of said environment.

### 4.3.1 The Environment

Snake is a relatively well-known game [18], but there are often slight differences between different versions. The following is explaining my particular game. Everything plays within a small grid, where a snake has to try and eat a fruit, which spawns and respawns at a random (non-occupied) location every time it gets eaten by the snake. Initially, the snake only consists of a head, but it grows by one body unit with every fruit eaten. At all times, there are four different actions available to the agent: up, down, right, or left. They all result in a movement by one unit. The body will follow the head's movement with a matching delay. The score is defined as the number of body units. Note that in this case, the score and the rewards received by the agent do not correlate in the same way as in the previous experiments.



Figure 24: Snake snapshot, score = 5

In Figure 24, we can see a snapshot of the game, where the snake consists of five body units, and thus, the score is equal to five. In this case, it plays in a seven by seven grid, where the maximum achievable score is 48, at which point the game would end due to winning. However, there are two more ways the game can end. Both result from failure. For one, if the head hits the border (i.e., tries to leave the grid), and second if the head hits its own body.

There are two more things to consider when designing the environment. How will the state observed by the agent look like? Which type of actions should be rewarded and which actions should be punished? In this experiment, the answers to these questions will be considered parameters. There will be many options to choose from, especially for the state representation, but more on that later.

### 4.3.2 The Agent

We will now look at the agent since he has more aspects to consider than those in the previous experiments. The agent is essentially made up of two main components, he has a Q-network, and he knows how to use the Q-network and optimize it based on the implemented algorithm. These two components give rise to the agent's parameters. Following is a list I curated with all the relevant parameters.

- Environment
  - State type
  - Rewards
- Agent
  - Q-Network, size and amount of hidden layers
  - Size of the agent's memory $\mathcal{D}$
  - Discount factor, $\gamma$
  - Learning rate, $\eta$
  - Exploration rate, $\varepsilon$
  - Number of weight updates per step, $n$

Now similar to the previous experiments, all one has to do is to try and find the best performing parameters through trial and error. This is easily said, but in reality, this is the part that often takes the longest, especially in experiments where the task at hand is rather complicated, which results in long training times until we see any sort of results from the agent. This is also dependent on the hardware being used.

### 4.3.3 State Types

We will now look at what I found to be the most interesting parameter, the state type. The question here is, how should the agent perceive his environment. There are two main ways of going about this. One is to give the agent a top-down view, i.e., what we humans see when we play snake. The second option is to embrace the 2D environment and imagine what one would see in it. For example, we could give the agent the distances from his head to the border in eight different directions, and we could also add binary vision for his body and the fruit. Binary vision means that going from his head in a specific direction, is there a fruit or not?

border distance     binary fruit     binary body

Figure 25: Snake state type, 2D vision

A simple example to illustrate and further understand the second concept is visible in Figure 25. There we can see a graphic for each object included in the vision. In the left one, we can see how the distances to the border are being measured. In the second and third ones, the binary vision is represented by the black and yellow lines. Yellow means that there is an object detected, and black means there is not. In total, this creates an input vector with 24 elements, where every line represents an element. The distances are simply real numbers. The binary vision lines are either a 1 if an object is detected or a 0 otherwise. Note that this is simply an example and this concept does not stipulate the need for eight directions or for the binary vision. It could also be only four directions for the border distances and maybe a distance measure for the fruit instead of a binary one.

We will now quickly look at the different options for the top-down representation for the agent's state. The first one that comes to mind is capturing the whole grid and feeding it to the agent, where different numbers represent the different objects on the grid. However, there are once again many ways to go about this. One that I found to work better than the former is only to show the agent a part of the grid centered at his head. This representation solves two problems with the previous one. First, the need to specifically represent the snake's head is no more since it is always at the center of the agent's visible sub-grid. Second, this state representation is independent of the grid size. This gives the advantage that an agent trained on a five by five grid can also play on a nine by nine or any other size. In Figure 26 one can see an illustration of the just discussed representation, where the agent would only be able to see the yellow highlighted sub-grid centered at the head.

Figure 26: Snake top-down view, 3x3 sub-grid centered at the head

I will now present to you the different state types that were used in the later following results. Note: Blue lines are distance measurements, black/yellow lines are binary vision, yellow highlighted areas represent a top-down view, cyan arrows are 2D vectors from the head to the fruit.



Figure 27: Snake state type #1



Figure 28: Snake state type #2

Figure 29: Snake state type #3



Figure 30: Snake state type #4



Figure 31: Snake state type #5



Figure 32: Snake state type #6

Figure 33: Snake state type #7



Figure 34: Snake state type #8

### 4.3.4 Parameter Optimization

I began by optimizing every parameter but the state type to compare the different state types with an optimized set of parameters in the end. Note that the "optimized set of parameters" is by no means optimal. It is merely what I found to work best within a limited time of trial and error testing. Also, I had to choose one state type to optimize the other parameters with and then operated under the assumption that the parameters work similarly well on the other state types. These are compromises I had to make to keep it realistic time-wise. Following is the set of optimized parameters, which I later used to compare the different state types.

- Rewards:
    - -50 for loosing
    - 50 for eating a fruit
    - 100 for winning
    - 0 otherwise
- Q-Network, 1 hidden layer with 128 neurons
- $|\mathcal{D}| = 10'000$
- $\gamma = 0.98$

- $\eta = 5\text{e-}5$
- $\varepsilon = (0.9998)^s$, $s$ increments with every step, starting at 0
- $n = 32$, weight updates per step

One further optimization I found to impact performance significantly is to let the agent start his training on a small grid and then work his way up to the final grid size, instead of immediately dropping him into the final grid size. This often causes an increase in performance, because on a small grid the agent is more likely to experience a positive reward (i.e., eating a fruit) than on a large grid. This is especially true in the beginning when the agent acts randomly.

### 4.3.5   Training Methodology

For each state type, two agents will be trained to check for the run-to-run variance. All 16 of them will be trained in the same manner except for the two with state type #4. The reason being is that an agent with said state type can only play in one specific grid size since we cannot change the input dimensions. Once we create his Q-network with a specific amount of input neurons, for example, in a 3x3 grid, we need nine input neurons, one for each cell. The agent is then bound to that specific grid size as we cannot change the input neurons of his Q-network. That would require initializing a new and untrained one. So all the agents except for the two with state type #4 are trained as follows.

- 200 episodes on a 3x3 grid
- 400 episodes on a 5x5 grid
- 800 episodes on a 7x7 grid
- 1600 episodes on a 9x9 grid

The two agents with state type #4 will spend all 3000 episodes on a 9x9 grid. As one can probably tell, the goal for all agents is to be able to play on a 9x9 grid.

### 4.3.6   Results

The following is an overview of the agent's performance during training.

Figure 35: Snake, training performances

From Figure 35 one can already get an impression of the performances. But to thoroughly test the agents, each of them played 1000 episodes in a testing mode. That means the exploration rate is zero, thus they always act greedily. The results from that testing and further information about the agents are presented in Table 2.

| agent ID | weight updates in millions | training duration in minutes | high score during training | average score during testing | high score during testing |
|---|---|---|---|---|---|
| 101.1 | 18.9 | 558 | 28 | 10.016 | 26 |
| 101.2 | 17.79 | 534 | 31 | 13.539 | 28 |
| 102.1 | 12.62 | 401 | 32 | 17.166 | 33 |
| 102.2 | 12.71 | 403 | 29 | 14.168 | 32 |
| 103.1 | 17.53 | 485 | 28 | 10.266 | 34 |
| 103.2 | 21.67 | 583 | 27 | 5.179 | 28 |
| 104.1 | 49.62 | 1055 | 4 | 0.417 | 4 |
| 104.2 | 55.07 | 1287 | 5 | 0.325 | 3 |
| 105.1 | 17.07 | 436 | 44 | 23.590 | 42 |
| 105.2 | 17.39 | 444 | 40 | 19.288 | 41 |
| 106.1 | 9.88 | 283 | 38 | 21.846 | 44 |
| 106.2 | 10.34 | 293 | 37 | 21.812 | 45 |
| 107.1 | 13.16 | 304 | 46 | 26.652 | 46 |
| 107.2 | 13.57 | 314 | 47 | 27.575 | 49 |
| 108.1 | 13.1 | 303 | 46 | 27.101 | 59 |
| 108.2 | 13.33 | 308 | 43 | 26.844 | 51 |

Table 2: Snake, final results

Especially visible in Figure 35, agent 104.1 and 104.2 performed significantly worse. That is partly because the state type did not allow for the previously discussed optimization strategy, where the agent starts on a small grid, which slowly grows to the final grid size. In a way this shows the immense importance of this strategy, but we cannot ignore the fact that the parameters besides the state type were optimized using this strategy. Also the state type itself may simply be a bad performing one. Therefore we cannot fully credit the lack of the optimization for the significantly worse performance.

Besides the failure of the state type #4 there are other notable results we can conclude from Figure 35 and Table 2. We can confidently say that the top-down view overall delivered better results than the 2D vision. We can also conclude that the addition of the vector between the snake's head and the fruit in state types #6 and #8 compared to #5 and #7 resulted in negligible improvements, if any. But, the most important question to answer with these results is the previously stated hypothesis. This is more complicated than anticipated

because we somehow have to define human-level performance in the snake game. To do so, I asked 15 people, including classmates and the teacher of a game programming class, to play my specific Snake version. Each of them played 30 episodes, and here are the results.

| overall average score (450 episodes) | best player's average score (30 episodes) | high score (within 450 episodes) |
|---|---|---|
| 1.1 | 3.17 | 19 |

Table 3: Human test results

From these results, we can conclude that even the worse state types, except for #4, managed to outperform the human test group by a significant margin, especially when we compare our on average best performing agent 107.2, with an average score of 27.575 to the human average of 1.1. With these results I can confirm my second hypothesis.

To further investigate the failure of #4, I was curious to see how the other state types would perform if they were trained on a 9x9 for 3000 episodes. Therefore I trained four more agents with state type #1, #2, #6, and #8. Following are the results presented in an identical manner.



Figure 36: Snake, training performances for the investigation of state type #4

44

| agent ID | weight updates in millions | training duration in minutes | high score during training | average score during testing | high score during testing |
|----------|---------------------------|------------------------------|----------------------------|------------------------------|---------------------------|
| 201 | 21.45 | 622 | 33 | 8.081 | 31 |
| 202 | 13.99 | 448 | 29 | 14.741 | 29 |
| 206 | 14.49 | 460 | 41 | 22.666 | 42 |
| 208 | 30.82 | 846 | 37 | 23.264 | 39 |

Table 4: Snake, final results for the investigation of state type #4

Figure 36 and Table 4 show how agents 201, 202, 206 performed similar to their counterparts (101.1, 101.2, 102.1, ...). However, with agent 208, we see different behavior, likely due to the large state type. It took him significantly longer to show any sign of learning. Nevertheless, he did manage to perform in the end, though his counterparts (108.1, 108.2) performed significantly better not only during training but also during testing. Nevertheless, we can now safely conclude that state type #4 is simply a bad choice, and its bad performance is not entirely due to the lack of the previously discussed optimization strategy.

# 5  Conclusion

Motivated by a single question: *How much can a Matura student with hardly any amount of prior knowledge in the field of reinforcement learning achieve in said field within the scope of a Matura Paper?* I sought to answer two hypotheses.

- One of the most basic solution methods of reinforcement learning, namely tabular Q-Learning, is adequate to solve the Cartpole environment, provided by OpenAI-Gym [1].
- Using Deep reinforcement learning, an agent can learn to play the game of Snake on a human-level performance.

As already discussed in Section 3.3, the Cartpole environment is said to be solved if the average score over the last 100 episodes is $\geq 195$. By Implementing the tabular Q-learning algorithm (Algorithm 1), the agent solved the environment in 164 episodes. In a paper by Swagat Kumar [5], he managed to exceed the 195 threshold in 300 episodes with the standard tabular Q-learning algorithm and in about 150 episodes with a deep Q-learning algorithm utilizing Prioritized Experience Replay (PER). Considering all of that, I found my results to be satisfactory and I was able to confirm my first hypothesis.

In Section 4.3, I managed to confirm my second hypothesis. There I implemented a deep Q-learning algorithm (Algorithm 2) on a self-made Snake game. I then began by optimizing every parameter but the state type to compare the different state types in the end. There, one was able to see the significance of the state representation and the sometimes drastic differences in performance that they caused. Finally, I managed to create an agent, which outperformed a test group of humans by a significant margin. The agent's (ID: 108.1) average score was 27.1, which translates to filling 35% of the grid. Its impressive high score was 59, which is equivalent to filling 74% of the grid. For comparison, I found a paper [15], in which the authors implemented a similar deep Q-learning algorithm. Their agent managed to achieve an average score of 9.04, which in their Snake game translates to filling 8% of the grid. Its high score was 17, which again translates to filling 14% of the grid. The most significant difference in their implementation, compared to mine, was that they took the raw pixels (240x240) as the state representation and combined that with convolutional layers in their neural network. To sum up I was once again very please with my results.

# 6  Reflection

The journey this work took me on was genuinely extraordinary. As I had just embarked upon this journey, I could have never imagined all the things I would learn. I started with no knowledge in the field of reinforcement learning nor in deep learning. Also, I had only just discovered programming for myself and had scarcely any amount of experience with Python. I had just finished my first small project, where I programmed the Snake game, which, in a way, inspired me to create an agent who could beat me in it. However, I frankly had no idea whether it was doable or not. I only had a goal but no clear path to it. Nevertheless, I think that it turned out great. I started small with a simple environment and a tabular Q-learning algorithm and worked my way up to implementing a deep Q-learning algorithm in my self-made Snake game.

On this journey, I did not only improve my Python skills by a significant margin. I also got to know LaTeX, which amazed me in many ways as a text writing tool. However, besides all the fantastic software and Python libraries I discovered, I am also glad that I took my time to understand the theory behind it all. Where the book by Richard Sutton and Andrew Barto [10] as well as the cost-free online available course by David Silver [9] played an immense and crucial role.

Finally, I would like to extend my sincere thanks to my supervisor Urs Sieber for his professional advice. I especially appreciate that I was able to work on my project independently.

## Declaration of Original Work

I confirm with my signature that this paper is entirely my own work and that any assistance given by others was restricted to advice and proof-reading. All sources employed in preparation of the paper and all quotations used are clearly cited and due acknowledgment is given for all help provided by others. I am aware of the definition of plagiarism in the Matura paper guidelines and that submission of work which has been plagiarised is a serious breach of the Matura regulations (Art. 1quarter of the Maturitätsprüfungsreglements des Gymnasiums).

Date and signature: ...........................................................................................................

# References

[1] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.

[2] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323, 2011.

[3] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[4] R. Koppula. Exploration vs. exploitation in reinforcement learning. `https://www.manifold.ai/exploration-vs-exploitation-in-reinforcement-learning`. [Online; accessed 5-November-2020].

[5] S. Kumar. Balancing a cartpole system with reinforcement learning–a tutorial. *arXiv preprint arXiv:2006.04938*, 2020.

[6] L.-J. Lin. Reinforcement learning for robots using neural networks. Technical report, Carnegie-Mellon Univ Pittsburgh PA School of Computer Science, 1993.

[7] Melissa and D. Store. Shape sorting cube. `https://www.amazon.co.uk/Melissa-Doug-Shape-Sorting-Cube/dp/B0731PF3TR`, 2020. [Online; accessed 09-September-2020].

[8] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

[9] D. Silver. Introduction to rl with david silver. `https://www.youtube.com/watch?v=2pWv7GOvuf0&list=PLqYmG7hTraZBiG_XpjnPrSNw-1XQaM_gB&ab_channel=DeepMind`, 2015. [Online; accessed 12-September-2020].

[10] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[11] C. Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.

[12] G. Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.

[13] C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.

[14] C. J. C. H. Watkins. Learning from delayed rewards. 1989.

[15] Z. Wei, D. Wang, M. Zhang, A.-H. Tan, C. Miao, and Y. Zhou. Autonomous agents in snake game via deep reinforcement learning. In *2018 IEEE International Conference on Agents (ICA)*, pages 20–25. IEEE, 2018.

[16] Wikipedia contributors. Markov decision process — Wikipedia, the free encyclopedia. `https://en.wikipedia.org/w/index.php?title=Markov_decision_process&oldid=975754122`, 2020. [Online; accessed 12-September-2020].

[17] Wikipedia contributors. Reinforcement learning — Wikipedia, the free encyclopedia. `https://en.wikipedia.org/w/index.php?title=Reinforcement_learning&oldid=976038269`, 2020. [Online; accessed 12-September-2020].

[18] Wikipedia contributors. Snake (video game genre) — Wikipedia, the free encyclopedia. `https://en.wikipedia.org/wiki/Snake_(video_game_genre)`, 2020. [Online; accessed 14-December-2020].

# List of Figures

# List of Tables

# List of Codes

# A   Source Code: Gridworld Environment

```python
1  ############################################################
2  # This code contains my self-made gridworld environment
3  ############################################################
4
5
6  import numpy as np
7  import drawSvg as draw
8
9  #inspired by the Gridworld example found in:
10 #Sutton, Richard S., and Andrew G. Barto. Reinforcement learning: An
       introduction. MIT press, 2018.
11
12 class Gridworld:
13     def __init__(self, m=7, n=14, goal_cords=[3,13], wall_cords_list=[]):
14         #initializing grid, player's position, inner walls,
15         #and the goal's position
16
17         self.m = m
18         self.n = n
19         self.n_actions = 4
20         self.game_matrix = np.full((m,n),fill_value=-1)
21         self.player_pos = [m//2, 0]
22         self.player_pos_temp = [m//2, 0]
23         self.wall_cords_list = wall_cords_list
24         self.goal_cords = goal_cords
25         if len(wall_cords_list) == 0:
26             for i in range(6):
27                 wall_cords_list.append([i,6])
28                 wall_cords_list.append([6-i,10])
29
30
31     def move_player(self, action):
32         #taking in an action and returning the player's theoratical position
33         self.player_pos_temp = self.player_pos.copy()
34         if action == 0:
35             self.player_pos_temp[0] -= 1
36         elif action == 1:
37             self.player_pos_temp[1] += 1
38         elif action == 2:
39             self.player_pos_temp[0] += 1
40         elif action == 3:
```

```python
41                 self.player_pos_temp[1] -= 1
42             else:
43                 print('invalid action')
44
45     def step(self, action):
46         #this function adjusts the environment based on the agent's action
47         #and returns: his new position, reward and a boolean to indicate
48         #whether the episode is done or not
49
50         self.move_player(action)
51
52         #checking whether the agent stepped into an inner wall
53         for cords in self.wall_cords_list:
54             if cords == self.player_pos_temp:
55                 reward = -100
56                 done = True
57                 self.player_pos = self.player_pos_temp
58                 return self.player_pos, reward, done
59
60         #checking whether the agent tried to leave the grid
61         if (self.player_pos_temp[0] >= self.m or
62             self.player_pos_temp[0] < 0 or
63             self.player_pos_temp[1] >= self.n or
64             self.player_pos_temp[1] < 0):
65             reward = -10
66             done = False
67             return self.player_pos, reward, done
68
69         #checking whether the agent reached the goal
70         elif self.player_pos_temp == self.goal_cords:
71             reward = -1
72             done = True
73             self.player_pos = self.player_pos_temp
74             return self.player_pos, reward, done
75
76         #if none of the above check then it was a "standard" action
77         else:
78             reward = -1
79             done = False
80             self.player_pos = self.player_pos_temp
81             return self.player_pos, reward, done
82
83     def reset(self):
```

```
84        #reseting the agent's position (the environment)
85        self.player_pos = [self.m//2, 0]
86        self.player_pos_temp = [self.m//2, 0]
87        return self.player_pos
88
89    def print_matrix(self):
90        #for developing purposes
91        temp = np.full((self.m,self.n), fill_value=-1)
92        temp[tuple(self.goal_cords)] = 0
93        for cords in self.wall_cords_list:
94            temp[tuple(cords)] = -100
95        temp[tuple(self.player_pos)] = 1
96        print(temp)
97
98    def print_matrix_symbols(self):
99        #to quickly check how it is looking and where the agent is located
100       #9 = goal, -9 = inner wall, 1 = agent
101       temp = np.zeros((self.m,self.n), dtype=int)
102       temp[tuple(self.goal_cords)] =  9
103       for cords in self.wall_cords_list:
104           temp[tuple(cords)] = -9
105       temp[tuple(self.player_pos)] = 1
106       print(temp)
107
108   def draw(self, path, save=True):
109       #saves a .svg file of the current environment state
110
111       d = draw.Drawing(142,72, origin=(-1,-1)) #creating canvas
112
113       #drawing inner grid
114       for i in range(1,self.m):
115           d.append(draw.Line(0,10*i,140,10*i,
116                       stroke='black', stroke_width=0.5))
117       for i in range(1,self.n):
118           d.append(draw.Line(10*i,0,10*i,70,
119                       stroke='black', stroke_width=0.5))
120
121       #drawing outer walls
122       d.append(draw.Line(0,0,140,0, stroke='black', stroke_width=1))
123       d.append(draw.Line(0,10*self.m,140,10*self.m,
124                       stroke='black', stroke_width=1))
125
126       d.append(draw.Line(0,0,0,70, stroke='black', stroke_width=1))
```

```
127        d.append(draw.Line(10*self.n,0,10*self.n,70,
128                           stroke='black', stroke_width=1))
129
130        #drawing inner walls
131        for cords in self.wall_cords_list:
132            d.append(draw.Line(10*(cords[1]),10*(self.m-1-cords[0]),
133                               10*(cords[1]+1),10*(self.m-1-cords[0]+1),
134                                                  stroke='black'))
135            d.append(draw.Line(10*(cords[1]),10*(self.m-1-cords[0]+1),
136                               10*(cords[1]+1),10*(self.m-1-cords[0]),
137                                                  stroke='black'))
138
139        #drawing player and goal
140        d.append(draw.Circle(self.player_pos[1]*10+5,
141                             (self.m-1-self.player_pos[0])*10+5,3.5,
142                             stroke='blue', fill='white', stroke_width=1))
143        d.append(draw.Rectangle(self.goal_cords[1]*10+.75,
144                             (self.m-1-self.goal_cords[0])*10+.75, 8.25, 8.5,
145                             stroke='green', fill='green'))
146        if save:
147            d.saveSvg(path)
148        else:
149            return d
150
151    def draw_greedy_policy(self, q_table, path):
152        d = self.draw('', save=False)
153        arrow = draw.Marker(-0.1, -0.5, 0.9, 0.5, scale=4, orient='auto')
154        arrow.append(draw.Lines(-0.1, -0.5, -0.1, 0.5, 0.9, 0,
155                                fill='red', close=True))
156
157
158        states = [(i,j) for i in range(self.m) for j in range(self.n)]
159
160        for s in states:
161            m = np.argmax(q_table[s])
162            optimal_actions = []
163
164            for i,v in enumerate(q_table[s]):
165                if v == q_table[s][m]:
166                    optimal_actions.append(i)
167
168            if s == (3,0):
169                print(optimal_actions)
```

```
170
171              for o_a in optimal_actions:
172                  if o_a == 0:
173                      d.append(draw.Line(10*(s[1])+5,10*(self.m-1-s[0])+5,
174                                         10*(s[1])+5, 10*(self.m-1-s[0])+7,
175                                         stroke='black', stroke_width=0.5,
176                                            marker_end=arrow))
177                  elif o_a == 1:
178                      d.append(draw.Line(10*(s[1])+5,10*(self.m-1-s[0])+5,
179                                         10*(s[1])+7, 10*(self.m-1-s[0])+5,
180                                         stroke='black', stroke_width=0.5,
181                                            marker_end=arrow))
182                  elif o_a == 2:
183                      d.append(draw.Line(10*(s[1])+5,10*(self.m-1-s[0])+5,
184                                         10*(s[1])+5, 10*(self.m-1-s[0])+3,
185                                         stroke='black', stroke_width=0.5,
186                                            marker_end=arrow))
187                  elif o_a == 3:
188                      d.append(draw.Line(10*(s[1])+5,10*(self.m-1-s[0])+5,
189                                         10*(s[1])+3, 10*(self.m-1-s[0])+5,
190                                         stroke='black', stroke_width=0.5,
191                                            marker_end=arrow))
192          d.saveSvg(path)
193
194
195
196
197 if __name__ == '__main__':
198
199     ###
200     #demonstration of how to use the draw_greedy_policy function
201     ###
202
203     #initialize environment
204     env = Gridworld()
205
206     #load a q-table
207     q_table = np.load('q_table.npy', allow_pickle=True)
208
209     #pass the q-table to the function
210     #and name the file
211     env.draw_greedy_policy(q_table, 'gridworld_q0=-90.svg')
```

Code 1: Source Code for the Gridworld Environment

# B  Source Code: Gridworld Solver

```
1  ################################################################
2  # This code implements the tabular Q-learning algorithm in a
3  # self-made gridworl environment (described in Section 3.2)
4  ################################################################
5
6
7  import numpy as np
8  import matplotlib.pyplot as plt
9  from gridworld_env import Gridworld
10 import math
11
12 #partly inspired by:
13 #https://pythonprogramming.net/q-learning-reinforcement-learning-python-
       tutorial/
14
15 class GridworldSolver:
16     def __init__(self, q0=0, gamma=1, epsilon=0.1, alpha=0.1):
17         #initializing the gridworld, the parameters and other
18         #later useful variables
19         self.env = Gridworld()
20
21         self.n_actions = self.env.n_actions
22         self.q0 = q0
23
24         self.gamma = gamma
25         self.epsilon = epsilon
26         self.alpha = alpha
27
28         self.optimal_s0 = [-27, -25, -25, -35]
29
30     def init_q_table(self):
31         #initializing q-table with user-inputed value
32         self.q_table = np.full((((self.env.m, self.env.n, self.n_actions,)),
33                                             fill_value=self.q0, dtype=float)
34
35         #setting q-value of terminal states to 0
36         self.q_table[tuple(self.env.goal_cords)] = [0 for _ in range(self.
    n_actions)]
37         for cords in self.env.wall_cords_list:
38             self.q_table[tuple(cords)] = [0 for _ in range(self.n_actions)]
39
```

60

```python
40
41     def update_q_value(self, state, new_state, action, reward, done, alpha):
42         #function to update the q-table based on observations
43         self.q_table[state][action] += (alpha * (reward + self.gamma *
44                                          np.max(self.q_table[new_state]) -
45                                          self.q_table[state][action]))
46
47     def one_episode(self, epsilon, alpha):
48         #reseting the environment
49         done = False
50         state = tuple(self.env.reset())
51         score = 0
52         while not done:
53             #choosing an action based on e-greedy action selection
54             action = np.argmax(self.q_table[state])
55             if np.random.uniform() < epsilon:
56                 action = np.random.randint(0,self.n_actions)
57
58
59
60             #taking the selected action and observing its consequences
61             new_state, reward, done = self.env.step(action)
62             new_state = tuple(new_state)
63
64             #keeping track of the reward
65             score += reward
66
67             #updating the q-table
68             self.update_q_value(state, new_state, action, reward, done, alpha)
69
70             state = new_state
71
72         return score
73
74     def n_episodes(self, n):
75         #initializing q-table, parameters, and score_history
76         #to keep track of the agent's performance
77         self.init_q_table()
78
79         score_history = []
80         solved_after = -1
81         epsilon = self.epsilon
82         alpha = self.alpha
```

```
 83
 84        #playing n episodes and keeping track of the agent's performance
 85        for e in range(n):
 86
 87            score = self.one_episode(epsilon, alpha)
 88            score_history.append(score)
 89
 90            ###
 91            #uncomment the following, if one would like to save the q-table
 92            #at the point were it was "optimal" for the first time
 93            #"optimal" as described in Section 3.2
 94            ###
 95
 96            #if (list(self.q_table[(3,0)]) == self.optimal_s0
 97            #                          and solved_after == -1):
 98            #    np.save('q_table.npy', self.q_table)
 99            #    solved_after = e
100
101        return score_history, solved_after
102
103    def N_n_episodes(self, N, n):
104        #this is to test different sets of parameters,
105        #it takes the average of multiple runs
106        s = []
107        solved_average = []
108
109        for _ in range(N):
110            l = self.n_episodes(n)
111            s.append(l[0])
112            solved_average.append(l[1])
113
114        return list(np.average(s, axis=0)), np.average(solved_average)
115
116    def plot_scores(self, score_histories, labels, solved_after, title='plot')
     :
117        #plots the agent's score of the individual episodes
118        #the episode on which he solved
119
120        plt.rcParams.update({'font.size': 15})
121        for i,score_history in enumerate(score_histories):
122            plt.plot(score_history, label=f'{labels[0]}: {labels[1][i]}',
123                    color=f'C{i}')
124            if solved_after[i] != -1:
```

```
125              plt.plot([solved_after[i],solved_after[i]],[-50,0],
126                      color=f'C{i}', linestyle='--', linewidth=3,
127                      label=f'converged after {solved_after[i]} episodes')
128
129        plt.plot([-25 for _ in range(len(score_histories[0]))], color='k',
130                                        label='optimal return, -25')
131        plt.xlabel('episode')
132        plt.ylabel('score')
133        plt.legend(loc='lower right', fontsize='x-large')
134        plt.title(title)
135        plt.show()
136
137
138    def __str__(self):
139        #to check the currently used parameters
140        return (f'epsilon={self.epsilon}, alpha={self.alpha}, '+
141                f'q0={self.q0}, gamma={self.gamma}')
```

Code 2: Source Code for the Gridworld Solver

# C    Source Code: Gridworld Main Script

```python
############################################################
# This code is a demonstration how to use the GridworldSolver
# (utilizing multiprocessing to speed it up)
# To execute, please place gridworld_env.py, gridworld_solver.py
# and gridworld_solver.py in the same directory
############################################################


import gridworld_solver
import multiprocessing as mp
import time
import numpy as np

#here one can change the parameters other than the one being adjusted
#every parameter other than the one being adjusted must have a default value
def one_process(epsilon, q0=0, gamma=1, alpha=.5, N=100, n=20000):
    solver = gridworld_solver.GridworldSolver(q0, gamma, epsilon, alpha)
    return solver.N_n_episodes(N,n)

#this is needed in order for the multiprocessing to work
if __name__ == '__main__':

    #keeping track of important metrics
    s0 = time.time()
    score_histories = []
    solved_after = []

    #at index 0 is the parameter that is being adjusted as str
    #at index 1 is a list of different values to be tested for said parameter
    parameters = ['epsilon', [0.1,0.3,0.5,0.7]]

    #using multiprocessing to
    p = mp.Pool(mp.cpu_count())
    l = p.map(one_process, parameters[1])
    for game in l:
        score_histories.append(game[0])
        solved_after.append(game[1])

    #saving the data
    np.save(f'data_{parameters[0]}.npy', np.array([score_histories,
                                solved_after, parameters], dtype=object))
```

```
42
43     plotter = gridworld_solver.GridworldSolver()
44     #the tile has to be adjusted manually
45     plotter.plot_scores(score_histories, parameters, solved_after,
46                         title='epsilon=_, alpha=0.5, q0=0, gamma=1')
47
48     #to see how long it took
49     print(f'calculations took {round((time.time()-s0)/60)} minutes')
50
51 ###
52 #if one does not want to compare different values for a parameter
53 #and would simply like to try one specific set of parameters
54 #here is how to do it
55 ###
56
57 #set parameters and create solver object
58 solver = gridworld_solver.GridworldSolver(epsilon=0.3,alpha=0.9, q0=-90)
59 #play a specific amount of episodes
60 s, solved = solver.n_episodes(20000)
61 #save the q-table (optional)
62 np.save('q_table.npy', solver.q_table)
```

Code 3: Source Code for the Gridworld Main Script

## D   Source Code: Cartpole Solver

```python
#################################################################
# This code implements the tabular Q-learning algorithm in the
# Cartpole environment provided by OpenAI-Gym
#################################################################

#partly inspired by:
#https://pythonprogramming.net/q-learning-reinforcement-learning-python-
    tutorial/


import numpy as np
import matplotlib.pyplot as plt
import gym
import math
import multiprocessing as mp
import sys

class CartpoleSolver:
    def __init__(self, obs_high=[2.4, 3.5, 0.21, 1.2],
                 obs_low=[-2.4, -3.5, -0.21, -1.2],
                 obs_chunks=[10, 10, 10, 10],
                 q0=0, gamma=1, epsilon=1, alpha=1, epsilon_min=0.1,
                 alpha_min=0.1, decrease_methode=[0,0], S_T_reward=-1000):

        #initializing variables
        self.n_actions = 2
        self.obs_high = np.array(obs_high)
        self.obs_low = np.array(obs_low)
        self.obs_chunks = np.array(obs_chunks)
        self.obs_chunks_delta = (self.obs_high - self.obs_low) / obs_chunks
        self.q0 = q0
        self.gamma = gamma
        self.epsilon = epsilon
        self.epsilon_min = epsilon_min
        self.alpha = alpha
        self.alpha_min = alpha_min
        self.state_list = []
        self.decrease_methode = decrease_methode
        self.S_T_reward = S_T_reward

        #initializing environment
```

```python
41          self.env = gym.make('CartPole-v0')

42

43      def discretize_state(self, state):
44          #turning a continuous state into a descrete one
45          discrete_state = list(((state - self.obs_low) /
46                          self.obs_chunks_delta).astype(np.int))
47          discrete_state = [min(self.obs_chunks[i]-1,
48                      discrete_state[i]) for i in range(4)]
49          return tuple(discrete_state)

50

51      def update_q_value(self, q_table, state, new_state,
52                      action, reward, done, alpha):
53          #updating the Q-table
54          if done:
55              #if it is done then the q_table[next_state]
56              #values should all be zero, hence that part is missing
57              q_table[state][action] += (alpha *
58                                  (reward - q_table[state][action]))
59          else:
60              q_table[state][action] += (alpha *
61                                  (reward +
62                                  self.gamma*np.max(q_table[new_state])
    -
63                                  q_table[state][action]))
64          return q_table

65

66      def decrease(self, v, v_min, episode, dm):
67          #function to decrease a parameter over time
68          #0-3 are different decreasing methods to choose from
69          if dm == 0:
70              return np.maximum(v_min,
71                          np.minimum(1., 1.0 - math.log10((episode + 1) / 25)))
72          elif dm == 1:
73              return np.maximum(v_min, v*(0.99)**episode)
74          elif dm == 2:
75              return np.maximum(v_min, 1.0 / np.sqrt(episode+1))
76          elif dm == 3:
77              return 1/(episode+1)

78

79      def one_episode(self, epsilon, alpha, render=False):
80          #reseting the environment
81          done = False
82          state = self.discretize_state(self.env.reset())
```

```
83          score = 0
84          while not done:
85              if render:
86                  self.env.render()
87
88              #choosing an action based on e-greedy action selection
89              action = np.argmax(self.q_table[state])
90              if np.random.uniform() < epsilon:
91                  action = self.env.action_space.sample()
92
93              #taking the selected action and observing its consequences
94              new_state, reward, done, _ = self.env.step(action)
95              self.state_list.append(new_state)
96              new_state = self.discretize_state(new_state)
97
98              #keeping track of the reward
99              score += reward
100
101             #adding extra punishment if the agent failed
102             if done and score < 200:
103                 reward = -1000
104
105             #updating the q-table
106             self.update_q_value(self.q_table, state, new_state,
107                                 action, reward, done, alpha)
108
109             state = new_state
110
111         return score
112
113     def n_episodes(self, n, render=False):
114         #initializing q-table, parameters, and score_history
115         #to keep track of the agent's performance
116         self.q_table = np.full((tuple(self.obs_chunks)+(self.n_actions,)),
117                                fill_value=self.q0, dtype=float)
118         score_history = []
119         epsilon = self.epsilon
120         alpha = self.alpha
121
122         #playing n episodes and keeping track of the agent's performance
123         for e in range(n):
124
125             score = self.one_episode(epsilon, alpha, render=render)
```

```
126            score_history.append(score)
127
128            #this would stop the agent and save his q_table
129            #as soon as he reaches the 195 threshold
130            #if e > 100:
131            #    if np.average(score_history[-100:]) >= 195:
132            #          np.save('q_table.npy', self.q_table)
133            #          print(e)
134            #          sys.exit()
135
136
137            epsilon = self.decrease(epsilon, self.epsilon_min,
138                                    e, self.decrease_methode[0])
139            alpha = self.decrease(alpha, self.alpha_min,
140                                  e, self.decrease_methode[1])
141
142
143        return score_history
144
145    def N_n_episodes(self, N, n, render=False):
146        #this is to test different sets of parameters,
147        #it takes the average of multiple runs
148
149        s = []
150
151        for _ in range(N):
152            s.append(self.n_episodes(n))
153
154        return list(np.average(s, axis=0))
155
156    def plot_runing_average(self, score_histories, labels, title='plot'):
157        #plots the average score from the last 100 episodes at each episodes
158        for score_history,label in zip(score_histories,labels[1]):
159            y = [np.average(score_history[max(0,i-100):i])
160                    for i in range(len(score_history))]
161            plt.plot(y, label=f'{labels[0]}: {label}',
162                color=f'C{score_histories.index(score_history)}')
163            plt.plot(y)
164
165        plt.plot([195 for _ in range(len(score_histories[0]))], color='k',
166                    label='195 threshold', linestyle=':')
167
168        plt.legend(loc='lower right', fontsize='x-large')
```

```
169         plt.xlabel('episode')
170         plt.ylabel('average return over last 100 episodes')
171         plt.title(title)
172         plt.show()
173
174     def plot_scores(self, score_histories, labels, title='plot'):
175         #plot the score of each episode
176         for score_history, label in zip(score_histories, labels[1]):
177             plt.plot(score_history, label=f'{labels[0]}: {label}',
178                     color=f'C{score_histories.index(score_history)}')
179         plt.legend()
180         plt.title(title)
181         plt.show()
182
183     def __str__(self):
184         #to check the currently used parameters
185         return (f'obs_high={self.obs_high}, obs_low={self.obs_low}, ' +
186                 f'obs_chunks={self.obs_chunks}, q0={self.q0}, ' +
187                 f'gamma={self.gamma}, epsilon={self.epsilon}, ' +
188                 f'alpha={self.alpha}, epsilon_min={self.epsilon_min}, ' +
189                 f'alpha_min={self.alpha_min}, S_T_reward={self.S_T_reward}, '+
190                 f'decrease_methode={self.decrease_methode}')
```

Code 4: Source Code for the Cartpole Solver

# E   Source Code: Cartpole Main Script

```python
####################################################################
# This code is a demonstration how to use the CartpoleSolver
# To execute, please place cartpole_solver.py and
# cartpole_main.py in the same directory
####################################################################


import cartpole_solver
import numpy as np

#to keep track of the scores
score_histories = []

#at index 0 is the parameter that is being adjusted as str
#at index 1 is a list of different values to be tested for said parameter
parameters_to_be_tried = ['alpha',[0.1,0.1,0.1,0.1,0.1]]

#iterates through the values to be tried
for p in parameters_to_be_tried[1]:
    solver = cartpole_solver.CartpoleSolver(
                        obs_high=[2.4, 2, 0.21, 1.8],
                        obs_low=[-2.4, -2, -0.21, -1.8],
                        obs_chunks=[1, 1, 12, 12], q0=0,
                        gamma=1, epsilon=1, alpha=p,
                        epsilon_min=0.01, alpha_min=p,
                        decrease_methode=[3,0], S_T_reward=-2000)
    #(amount of agents, amount of episodes) to be tested on a value
    l = solver.N_n_episodes(1,500)
    score_histories.append(l)
    print(solver.__str__())

#setting the title for the plot
title = f'adjusting {parameters_to_be_tried[0]}'

#saving the data (optional)
#so it can be laoded later to plot it again
np.save(f'data_{parameters_to_be_tried[0]}.npy',
        np.array([score_histories, parameters_to_be_tried, title],
                                            dtype=object))

#plotting the results
```

```
42 solver.plot_runing_average(score_histories,
43                                 parameters_to_be_tried, title=title)
44
45 ###
46 #example of loading data and plotting
47 ###
48
49 #initializing the CartpoleSolver object containing plotter function
50 plotter = cartpole_solver.CartpoleSolver()
51
52 #laoding the data
53 data = np.load('data_[name].npy', allow_pickle=True)
54
55 #setting the title and plotting the data
56 title = f'adjusting {data[1][0]}'
57 plotter.plot_runing_average(data[0],data[1],title)
58
59
60 ###
61 #alternatively one can test a single specific set of parameters
62 ###
63
64 #initializing solver with specific parameters
65 solver = cartpole_solver.CartpoleSolver(
66                 obs_high=[2.4, 2, 0.21, 1.8],
67                 obs_low=[-2.4, -2, -0.21, -1.8],
68                 obs_chunks=[1, 1, 12, 12], q0=0,
69                 gamma=1, epsilon=1, alpha=0.9,
70                 epsilon_min=0.01, alpha_min=0.9,
71                 decrease_methode=[3,0], S_T_reward=-2000)
72
73 #either by a single agent (n = amounts of episodes)
74 score_history = solver.n_episodes(n=1000,render=False)
75
76 #or by many agents (N = amount of agents)
77 score_histories = solver.N_n_episodes(N=10,n=1000,render=False)
```

Code 5: Source Code for the Cartpole Main Script

# F  Source Code: Snake Trainer

```
1  ################################################################
2  # This code comibines the Snake environment and the DQL agent
3  # It combines them in a python object, which can be used to
4  # train an agent with freely adjustable parameters
5  ################################################################
6
7
8  import snake_agent
9  import snake_env
10 import os
11 import torch
12 import time
13 import sys
14 import numpy as np
15
16 class SnakeTrainer:
17     def __init__(self, agent_id, updates_per_step, gamma, lr,
18                  epsilon, epsilon_decay, epsilon_min, statetype,
19                  nr_hiddenlayers, nr_neurons, max_mem_length,
20                  gridlength, r_fruit, r_collision, r_step,
21                  lr_decay=False, lr_milestones=None, lr_gamma=None,
22                  render_bool=False):
23
24         #selecting device, preferably the gpu
25         self.device = torch.device('cuda' if torch.cuda.is_available()
26                                             else 'cpu')
27
28         #paths for data collection
29         self.agent_id = agent_id
30         self.path = (os.path.dirname(os.path.abspath(__file__))
31                      + f'/results/{agent_id}/')
32
33         self.modelpath = self.path + 'checkpoint.tar'
34         self.settingspath = self.path + 'settings.txt'
35         self.datapath = self.path + 'data.tar'
36         self.memorypath = self.path + 'memory.tar'
37         self.replaypath = self.path + 'replay.npz'
38
39         #initializing important variables
40         self.updates_per_step = updates_per_step
41         self.render_bool = render_bool
```

```python
42
43          #initializing the environment
44          self.env = snake_env.Environment(gridlength=gridlength,
45              snakex=gridlength//2, snakey=gridlength//2,
46              fruitspawnset=False, r_fruit = r_fruit,
47              r_collision = r_collision, r_step = r_step,
48              statetype=statetype)
49
50          self.statetype = statetype
51          self.rewards_fcs = [r_fruit,r_collision,r_step]
52
53          obs_space = len(self.env.env_reset())
54
55          #initializing the agent
56          self.agent = snake_agent.Agent(gamma=gamma, lr=lr, epsilon=epsilon,
57              epsilon_decay=epsilon_decay, epsilon_min=epsilon_min,
58              obs_space=obs_space, act_space=4, device=self.device,
59              updates_per_step=self.updates_per_step, lr_decay=lr_decay,
60              lr_milestones=lr_milestones, lr_gamma=lr_gamma,
61              max_mem_length=max_mem_length,
62              nr_hiddenlayers=nr_hiddenlayers, nr_neurons=nr_neurons)
63
64          #initializing data collection
65          if not os.path.isdir(self.path):
66              self.init_data()
67          else:
68              self.load_agent()
69
70
71      def init_data(self):
72          #creates a directory and a first textfile with agent's specs
73          os.makedirs(self.path)
74          textfile = open(self.settingspath, 'a')
75          textfile.write(f"agent_id: {self.agent_id}"
76                  +f"device: {self.device} \nDQN: \n{self.agent.dqn}\n"
77                  +f"\noptimizer: \n{self.agent.optimizer}\n"
78                  +f"\ngamma: {self.agent.gamma} \nepsilon_decay & _min: "
79                  +f"{self.agent.epsilon_decay}, {self.agent.epsilon_min}"
80                  +f"\nstatetype: {self.env.statetype} \nupdates_per_step: "
81                  +f"{self.agent.updates_per_step} \nmax_mem_length: "
82                  +f"{self.agent.max_mem_length} \nrewards(f,c,s): "
83                  +f"{self.env.r_fruit}, {self.env.r_collision}, "
84                  +f"{self.env.r_step} \nlr_decay: {self.agent.lr_decay}")
```

```python
85          if self.agent.lr_decay:
86              textfile.write(f"lr_scheduler: \n{self.agent.scheduler}\n"
87                          +f"lr_milestones: {self.agent.lr_milestones} \nlr_gamma: "
88                          +f"{self.agent.lr_gamma}")
89          textfile.close()
90
91          #since its the first time this agent palys
92          # new data collection lists are created
93          self.episode_nr = 0
94          self.epsilon_history = []
95          self.stime_history = []
96          self.q_history = []
97          self.replay_memory = []
98          self.replay_memorytemp = []
99          self.updates = 0
100         self.main_time = 0
101
102     def load_agent(self):
103         #laods all that is necessary
104
105         #loading model
106         checkpoint = torch.load(self.modelpath)
107         self.agent.dqn.load_state_dict(checkpoint['model_state_dict'])
108         self.agent.optimizer.load_state_dict(
109                                 checkpoint['optimizer_state_dict'])
110         if self.agent.lr_decay:
111             self.agent.scheduler.load_state_dict(checkpoint['lr_scheduler'])
112         print("model was loaded")
113
114         #loading data
115         checkpoint = torch.load(self.datapath)
116         self.epsilon_history = checkpoint['epsilon_history']
117         self.stime_history = checkpoint['sruvival_time_history']
118         self.q_history = checkpoint['q_history']
119         self.env.scorelist = checkpoint['scorelist']
120         self.episode_nr = checkpoint['episode_nr']
121         self.agent.epsilon = checkpoint['epsilon']
122         self.updates = checkpoint['updates']
123         self.main_time = checkpoint['main_time']
124         print("data was loaded")
125
126         #loading agent's memory
127         self.agent.memory = torch.load(self.memorypath)
```

```python
128         print("memory was loaded")
129
130         #loading replay memory
131         self.replay_memorytemp = []
132         if os.path.exists(self.replaypath):
133             self.replay_memory = np.load(self.replaypath,
134                             allow_pickle=True)['replay'].tolist()
135         print("replay memory was loaded")
136
137     def save_model(self):
138         #saves the model
139         if self.agent.lr_decay:
140             torch.save({
141                 'model_state_dict': self.agent.dqn.state_dict(),
142                 'optimizer_state_dict': self.agent.optimizer.state_dict(),
143                 'lr_scheduler': self.agent.scheduler.state_dict(),
144                 }, self.modelpath)
145         else:
146             torch.save({
147                 'model_state_dict': self.agent.dqn.state_dict(),
148                 'optimizer_state_dict': self.agent.optimizer.state_dict()
149                 }, self.modelpath)
150
151     def save_data(self):
152         #saves the data
153         torch.save({
154             'episode_nr': self.episode_nr,
155             'epsilon': self.agent.epsilon,
156             'epsilon_history': self.epsilon_history,
157             'sruvival_time_history': self.stime_history,
158             'q_history': self.q_history,
159             'scorelist': self.env.scorelist,
160             'updates': self.updates,
161             'main_time': self.main_time
162             }, self.datapath)
163
164     def save_memory(self):
165         #save the memory
166         torch.save(self.agent.memory, self.memorypath)
167
168     def save_replay(self):
169         np.savez_compressed(self.replaypath,
170                     replay=np.array(self.replay_memory, dtype=object))
```

```
171
172
173     def train_n_episodes(self, n, gridlength, epsilon_reset=False):
174         #loading new environment incase the gridlength changed
175         self.env = snake_env.Environment(gridlength=gridlength,
176             snakex=gridlength//2, snakey=gridlength//2,
177             fruitspawnset=False, r_fruit = self.rewards_fcs[0],
178             r_collision = self.rewards_fcs[1], r_step = self.rewards_fcs[2],
179             statetype=self.statetype)
180
181         if os.path.exists(self.datapath):
182          self.load_agent()
183
184         if epsilon_reset:
185             self.agent.epsilon = 1
186
187         #starting the main for loop to play n episodes
188         for e in range(n):
189             self.episode_nr += 1
190             self.env.env_reset()
191             current_body = []
192             for b in self.env.body.list:
193                 current_body.append(b)
194             self.replay_memorytemp.append((self.env.snake.x,
195                     self.env.snake.y, self.env.fruit.x,
196                     self.env.fruit.y, current_body))
197
198             main_time_temp = time.time()
199
200             for i in range(1500):
201                 #receiving the observation from the environment
202                 #choosing an action and take it
203                 #receiving reward
204                 state = self.env.get_state()
205                 state = state.to(self.device)
206                 action = self.agent.act(state)
207                 reward = self.env.env_step(action)
208
209                 #keeping track of the estimated total reward at s_0
210                 if i == 0:
211                     self.q_history.append(torch.max(
212                                     self.agent.dqn.forward(state)).item())
213
```

```python
214                if self.render_bool:
215                    self.env.render(60, (0,0,0))
216
217                #saving the frame for later replay
218                current_body = []
219                for b in self.env.body.list:
220                    current_body.append(b)
221                self.replay_memorytemp.append((self.env.snake.x,
222                            self.env.snake.y, self.env.fruit.x,
223                                self.env.fruit.y, current_body))
224
225
226                #checking if he collided or fille the whole grid
227                if reward == self.rewards_fcs[1]:
228                    done = True
229                elif reward == 100:
230                    #incase the agent ever manages to fill the whole grid
231                    done = True
232                else:
233                    done = False
234
235                #receiving new state
236                new_state = self.env.get_state()
237                new_state = new_state.to(self.device)
238
239                #saving the agent's experience
240                self.agent.savedata(state, action,
241                                    reward, new_state, done)
242
243                #as soon as there is enough data the training will start
244                if len(self.agent.memory) > self.updates_per_step:
245                    self.agent.train()
246                    self.updates += self.updates_per_step
247
248                #if he collided with either the wall or himself,
249                # the episode is over
250                if done:
251                    break
252
253            #keeping track of important metrics
254            self.epsilon_history.append(self.agent.epsilon)
255            self.env.env_store_score()
256            self.stime_history.append(i+1)
```

78

```python
257            print(f"episode: {self.episode_nr}"
258                +f" score: {self.env.scorelist[-1]}"
259                +f" surv_time: {i+1}"
260                +f" epsilon: {round(self.agent.epsilon,3)}"
261                +f" q: {round(self.q_history[-1])}")
262
263
264            #keeping track of time taken
265            self.main_time += time.time()-main_time_temp
266
267            self.replay_memory.append(self.replay_memorytemp)
268            self.replay_memorytemp = []
269
270            #every 50th episode everything is saved
271            if self.episode_nr%50 == 0:
272                print("saving...")
273                self.save_model()
274                self.save_data()
275                self.save_memory()
276                self.save_replay()
277                print("done saving!")
278
279    def test_n_episodes(self, n, gridlength, record=False, render=True):
280        #loading new environment incase the gridlength changed
281        #recort = True will save a .png of every frame
282        self.env = snake_env.Environment(gridlength=gridlength,
283            snakex=gridlength//2, snakey=gridlength//2,
284            fruitspawnset=False, r_fruit = self.rewards_fcs[0],
285            r_collision = self.rewards_fcs[1], r_step = self.rewards_fcs[2],
286            statetype=self.statetype)
287
288        if os.path.exists(self.datapath):
289            self.load_agent()
290
291        score_list = []
292
293        import pygame
294        clock = pygame.time.Clock()
295        for e in range(n):
296            self.env.env_reset()
297            rewards = []
298            for i in range(1500):
299                if render:
```

```python
                    clock.tick(10)
                #receiving the observation from the environment
                #choosing an action and take it
                #receiving reward
                state = self.env.get_state()
                state = state.to(self.device)
                action = self.agent.act_greedy(state)
                reward = self.env.env_step(action)
                rewards.append(reward)

                #keeping track of the estimated total reward at s_0
                if i == 0:
                    q0 = torch.max(self.agent.dqn.forward(state)).item()


                if render:
                    self.env.render(60, (0,0,0), record=record,
                    path=self.path + f"frame{i}.png")



                #checking if he collided or fille the whole grid
                if reward == self.rewards_fcs[1]:
                    done = True
                elif reward == 100:
                    done = True
                else:
                    done = False

                #receiving new state
                new_state = self.env.get_state()
                new_state = new_state.to(self.device)

                #if he collided with either the wall or himself,
                # the episode is over
                if done:
                    break

            real_q0 = 0
            for j in range(len(rewards)):
                real_q0 += rewards[j] * (self.agent.gamma**(j))

            print(f'score: {self.env.score}, surv_time: {i},'
                + f' est q0: {round(q0)}, real q0: {round(real_q0)}')
            score_list.append(self.env.score)
```

```
343        return score_list
344
345    def watch_replay(self, episodes, gridlengths):
346        #used to watch a replay of a specific training episode
347        import pygame
348        pygame.init()
349        clock = pygame.time.Clock()
350        for e,g in zip(episodes,gridlengths):
351            print(self.env.scorelist[e])
352            for snakex,snakey,fruitx,fruity,bodylist in self.replay_memory[e]:
353                clock.tick(10)
354                self.env.renderreplay(gridlength=g, size=45,
355                                snakex=snakex, snakey=snakey,
356                                fruitx=fruitx, fruity=fruity,
357                                    bodylist=bodylist)
```

Code 6: Source Code for the Snake Trainer

# G   Source Code: Snake Agent

```python
1  ##############################################################
2  # This code contains the Python object for the DQL agent.
3  # the agent is made up of a DQN and the implemented algorithm
4  ##############################################################
5
6  #partly inspired by:
7  #https://github.com/philtabor/Youtube-Code-Repository/blob/master/
       ReinforcementLearning/DeepQLearning/simple_dqn_torch_2020.py
8
9
10 import torch
11 import torch.nn as nn
12 import torch.nn.functional as F
13 from collections import deque
14 import random
15 import numpy as np
16 import time
17 import os
18
19
20 class DQN(nn.Module):
21     def __init__(self, obs_space, act_space, nr_hiddenlayers, nr_neurons):
22         super(DQN, self).__init__()
23         #modular deep Q-network
24         #amount and size of linear hidden layers are freely adjustable
25         self.nr_hiddenlayers = nr_hiddenlayers
26         nr_neurons.insert(0,obs_space)
27         nr_neurons.append(act_space)
28         self.nr_neurons = nr_neurons
29
30         temp = []
31         for i in range(self.nr_hiddenlayers+1):
32             temp.append(nn.Linear(self.nr_neurons[i],self.nr_neurons[i+1]))
33
34         self.layers = nn.ModuleList(temp)
35
36
37     def forward(self, state):
38         #forward pass thorugh the network
39         x = state
40         for i in range(self.nr_hiddenlayers):
```

```python
41              x = F.relu(self.layers[i](x))

42

43          return self.layers[-1](x)

44

45

46  class Agent:
47      def __init__(self, gamma, lr, epsilon, epsilon_decay, epsilon_min,
48                   obs_space, act_space, device, updates_per_step,
49                   lr_decay = False, lr_milestones = None, lr_gamma = None,
50                   nr_hiddenlayers=1, nr_neurons=[128], max_mem_length=1000):
51          #deep q-learning agent
52          #initializing variables and constants
53          self.gamma = gamma
54          self.epsilon = epsilon
55          self.epsilon_decay = epsilon_decay
56          self.epsilon_min = epsilon_min
57          self.lr = lr
58          self.device = device
59          self.act_space = act_space
60          self.obs_space = obs_space
61          self.updates_per_step = updates_per_step
62          self.lr_decay = lr_decay
63          self.lr_milestones = lr_milestones
64          self.lr_gamma = lr_gamma
65          self.max_mem_length = max_mem_length

66

67          #initializing agent's memory
68          self.memory = deque(maxlen=self.max_mem_length)

69

70          #initializing agent's deep-Q-network
71          self.dqn = DQN(self.obs_space, self.act_space,
72                         nr_hiddenlayers, nr_neurons).to(self.device)
73          self.optimizer = torch.optim.Adam(self.dqn.parameters(), lr=lr)

74

75          #initializing learning rate scheduler, if lr decay is wanted
76          if self.lr_decay:
77              self.scheduler = torch.optim.lr_scheduler.MultiStepLR(
78                                        self.optimizer,
79                                        milestones=self.lr_milestones,
80                                        gamma=self.lr_gamma)

81

82

83      def savedata(self, state, action, reward,
```

```python
                                     new_state, gameovertemp):
        #used to append a transition to the agent's memory
        self.memory.append((state, action,
                            reward, new_state, gameovertemp))


    def act(self, state):
        #choosing action with epsilon-greedy strategy
        if np.random.rand() <= self.epsilon:
            return random.randrange(self.act_space)
        else:
            qvalues = self.dqn.forward(state)
            return torch.argmax(qvalues).item()


    def act_greedy(self, state):
        #choosing action completely greedy
        qvalues = self.dqn.forward(state)
        return torch.argmax(qvalues).item()


    def train(self):
        #sampling a specific number of datapoints from memory
        trainsample = random.sample(self.memory, self.updates_per_step)


        #iterates through every datapoint (transition) and
        # performs a weight update on the Q-network
        for state, action, reward, new_state, done in trainsample:
            if done:
                target = reward
            else:
                target = reward + (self.gamma
                        * torch.max(self.dqn.forward(new_state)).item())

            prediction = self.dqn.forward(state)[action].to(self.device)
            target = torch.tensor(target).to(self.device)
            self.optimizer.zero_grad()
            loss = F.mse_loss(target, prediction)
            loss.backward()
            self.optimizer.step()
            if self.lr_decay:
                self.scheduler.step()

        #adjustment of the exploration rate epsilon
        if self.epsilon > self.epsilon_min:
            self.epsilon = self.epsilon * self.epsilon_decay
```

```
127        else:
128            self.epsilon = self.epsilon_min
129
130    def get_lr(self):
131        #for development used to check current learning rate
132        for param_group in self.optimizer.param_groups:
133            return param_group['lr']
```

Code 7: Source Code for the Snake Agent

# H    Source Code: Snake Environment

```python
1  ################################################################
2  # This code contains my self-made Snake environment
3  # adjusted to be used as a Markov decision process
4  ################################################################
5
6
7  import random
8  import numpy as np
9  import torch
10 import drawSvg as draw
11
12 class Snake:
13     def __init__(self, spawnx, spawny):
14         #only considers the head and its movement
15         self.x = spawnx
16         self.y = spawny
17         self.changelist = []
18
19
20     def act(self, action):
21         #takes action and keep a record of it in the changelist
22         # for delayed body movement
23         if action == 0:
24             self.x += 1
25             self.changelist.append((1, 0))
26         if action == 1:
27             self.x -= 1
28             self.changelist.append((-1, 0))
29         if action == 2:
30             self.y += 1
31             self.changelist.append((0, 1))
32         if action == 3:
33             self.y -= 1
34             self.changelist.append((0, -1))
35
36 class Body:
37     def __init__(self, len):
38         #keeps track of the snake's body
39         self.len = len
40         self.list = []
41
```

```python
42     def grow(self, snake):
43         #function to grow the body
44         self.len += 1
45         self.list.append([snake.x,snake.y])
46
47     def update(self, changelist, snake):
48         #function to position each body part at its position
49         #uses the snake's head position and the changelist
50         # to calculate the position
51         for i in range(0,self.len):
52             x = snake.x
53             y = snake.y
54             for n in range(1,(i+2)):
55                 ch = changelist[-n]
56                 x = x-ch[0]
57                 y = y-ch[1]
58             self.list[i] = ([x,y])
59
60 class Fruit:
61     #used for the fruit
62     def __init__(self, gridlength, spawnset = False,
63                         spawnx = None , spawny = None):
64         #spawns the fruit at a either given or random position
65         if spawnset:
66             self.x = spawnx
67             self.y = spawny
68         else:
69             self.x = random.randint(0,(gridlength-1))
70             self.y = random.randint(0,(gridlength-1))
71
72 class Environment:
73     def __init__(self, gridlength, snakex, snakey, r_fruit,
74                 r_collision, r_step, fruitspawnset,
75                 fruitx = None, fruity = None, statetype = 1):
76
77         #initiliazing variables, snake, fruit and body
78         self.gridlength = gridlength
79         self.snakex = snakex
80         self.snakey = snakey
81         self.snakecolor = '#197416'
82         self.fruitcolor = '#E92323'
83         self.bodycolor = '#32B145'
84         self.background = '#485B57'
```

```python
85          self.fruitx = fruitx
86          self.fruity = fruity
87          self.fruitspawnset = fruitspawnset
88          self.snake = Snake(self.snakex, self.snakey)
89          self.fruit = Fruit(self.gridlength,
90                             spawnset=self.fruitspawnset,
91                             spawnx=self.fruitx, spawny=self.fruity)
92          self.body = Body(0)
93          self.score = 0
94          self.scorelist = []
95          self.gamewon = False
96          self.statetype = statetype
97
98          #adjustable rewards given to the agent
99          self.r_fruit = r_fruit
100         self.r_collision = r_collision
101         self.r_step = r_step
102
103     def collisionfruit(self):
104         #checks whether the snake has collided with the fruit
105         if self.snake.x == self.fruit.x and self.snake.y == self.fruit.y:
106             return True
107
108     def collisionborder(self):
109         #checks whether the snake has collided with the border
110         if (self.snake.x == self.gridlength or self.snake.x == -1 or
111             self.snake.y == self.gridlength or self.snake.y == -1):
112             return True
113
114     def collisionself(self):
115         #checks whether the snake has collided with its body
116         for b in self.body.list:
117             if self.snake.x == b[0] and self.snake.y == b[1]:
118                 return True
119
120     def check_fruit_spawn(self):
121         #makes sure the fruit can only spawn on a free cell
122         self.fruit = Fruit(self.gridlength)
123         for b in self.body.list:
124             if self.fruit.x == b[0] and self.fruit.y == b[1]:
125                 self.check_fruit_spawn()
126             else:
127                 pass
```

```
128            if self.fruit.x == self.snake.x and self.fruit.y == self.snake.y:
129                self.check_fruit_spawn()
130            else:
131                pass
132
133    def get_state(self):
134        #this function is needed,
135        # because there are different state types to choose from
136        if self.statetype == 1:
137            return self.get_state1()
138        elif self.statetype == 2:
139            return self.get_state2()
140        elif self.statetype == 3:
141            return self.get_state3()
142        elif self.statetype == 4:
143            return self.get_state4()
144        elif self.statetype == 5:
145            return self.get_state5()
146        elif self.statetype == 6:
147            return self.get_state6()
148        elif self.statetype == 7:
149            return self.get_state7()
150        elif self.statetype == 8:
151            return self.get_state8()
152
153    def get_state1(self):
154        #returns a 24-dimensional vector,
155        #consisting of the distances between the snake's head and:
156        #the border
157        #+ binary vision for the fruit and its body
158        #for each object 8 different directions get checked:
159        #if there is no object in that direction the distance is 0
160        #border distances to walls
161        #and diagonals
162
163        #walls, if it is 1 then one more step in that direction is lethal
164        b0 = float(self.gridlength-self.snake.x)
165        b2 = float(self.snake.y+1)
166        b4 = float(self.snake.x+1)
167        b6 = float(self.gridlength-self.snake.y)
168
169        b1 = False
170        b3 = False
```

```python
171         b5 = False
172         b7 = False
173         #diagonals
174         for i in range(self.gridlength):
175             if ((self.snake.x+i >= self.gridlength or
176                 self.snake.y+i >= self.gridlength) and
177                 (not b7)):
178                 b7 = float(np.sqrt(i**2+i**2))
179             if ((self.snake.x-i < 0 or
180                 self.snake.y-i < 0) and
181                 (not b3)):
182                 b3 = float(np.sqrt(i**2+i**2))
183             if ((self.snake.x+i >= self.gridlength or
184                 self.snake.y-i < 0) and
185                 (not b1)):
186                 b1 = float(np.sqrt(i**2+i**2))
187             if ((self.snake.x-i < 0 or
188                 self.snake.y+i >= self.gridlength) and
189                 (not b5)):
190                 b5 = float(np.sqrt(i**2+i**2))
191
192         f0 = 0.0
193         f1 = 0.0
194         f2 = 0.0
195         f3 = 0.0
196         f4 = 0.0
197         f5 = 0.0
198         f6 = 0.0
199         f7 = 0.0
200
201         s0 = 0.0
202         s1 = 0.0
203         s2 = 0.0
204         s3 = 0.0
205         s4 = 0.0
206         s5 = 0.0
207         s6 = 0.0
208         s7 = 0.0
209
210         if self.snake.y == self.fruit.y:
211             if self.snake.x < self.fruit.x:
212                 f0 = 1.0
213             else:
```

```
214                    f4 = 1.0
215            if self.snake.x == self.fruit.x:
216                if self.snake.y < self.fruit.y:
217                    f6 = 1.0
218                else:
219                    f2 = 1.0
220
221
222            for i in range(self.gridlength):
223                if (self.snake.x+i == self.fruit.x and
224                    self.snake.y+i == self.fruit.y):
225                    f7 = 1.0
226                if (self.snake.x-i == self.fruit.x and
227                    self.snake.y-i == self.fruit.y):
228                    f3 = 1.0
229                if (self.snake.x+i == self.fruit.x and
230                    self.snake.y-i == self.fruit.y):
231                    f1 = 1.0
232                if (self.snake.x-i == self.fruit.x and
233                    self.snake.y+i == self.fruit.y):
234                    f5 = 1.0
235
236            for b in self.body.list:
237                if self.snake.y == b[1]:
238                    if self.snake.x < b[0]:
239                        s0 = 1.0
240                    else:
241                        s4 = 1.0
242                if self.snake.x == b[0]:
243                    if self.snake.y < b[1]:
244                        s6 = 1.0
245                    else:
246                        s2 = 1.0
247                for i in range(self.gridlength):
248                    if (self.snake.x+i == b[0] and
249                        self.snake.y+i == b[1]):
250                        s7 = 1.0
251                    if (self.snake.x-i == b[0] and
252                        self.snake.y-i == b[1]):
253                        s3 = 1.0
254                    if (self.snake.x+i == b[0] and
255                        self.snake.y-i == b[1]):
256                        s1 = 1.0
```

```
257                  if (self.snake.x-i == b[0] and
258                      self.snake.y+i == b[1]):
259                      s5 = 1.0
260
261          return torch.tensor([b0, b1, b2, b3, b4, b5, b6, b7,
262                               f0, f1, f2, f3, f4, f5, f6, f7,
263                               s0, s1, s2, s3, s4, s5, s6, s7])
264
265      def get_state2(self):
266          #four directions for distances to border
267          #four directions for binary vision to body
268          #relative coordinates to fruit
269          #up, down, right and left
270
271          #border
272          b0 = float(self.gridlength-self.snake.x)
273          b1 = float(self.snake.y+1)
274          b2 = float(self.snake.x+1)
275          b3 = float(self.gridlength-self.snake.y)
276
277          #own body
278          s0 = 0.0
279          s1 = 0.0
280          s2 = 0.0
281          s3 = 0.0
282
283          for b in self.body.list:
284              if self.snake.x == b[0]:
285                  if self.snake.y >= b[1]:
286                      s1 = 1.0
287                  if self.snake.y <= b[1]:
288                      s3 = 1.0
289              if self.snake.y == b[1]:
290                  if self.snake.x <= b[0]:
291                      s0 = 1.0
292                  if self.snake.x >= b[0]:
293                      s2 = 1.0
294
295          #fruit
296          fx = float(self.fruit.x - self.snake.x)
297          fy = float(self.fruit.y - self.snake.y)
298          return torch.tensor([b0, b1, b2, b3,
299                               s0, s1, s2, s3,
```

```python
                                  fx, fy])

    def get_state3(self):
        #for border and body  in 4 directions
        #border with distances and body with binary vision
        #for the fruit a binary vision in 8 directions

        #border
        b0 = float(self.gridlength-self.snake.x)
        b1 = float(self.snake.y+1)
        b2 = float(self.snake.x+1)
        b3 = float(self.gridlength-self.snake.y)

        #body
        s0 = 0.0
        s1 = 0.0
        s2 = 0.0
        s3 = 0.0

        for b in self.body.list:
            if self.snake.x == b[0]:
                if self.snake.y >= b[1]:
                    s1 = 1.0
                if self.snake.y <= b[1]:
                    s3 = 1.0
            if self.snake.y == b[1]:
                if self.snake.x <= b[0]:
                    s0 = 1.0
                if self.snake.x >= b[0]:
                    s2 = 1.0

        #fruit
        f0 = 0.0
        f1 = 0.0
        f2 = 0.0
        f3 = 0.0
        f4 = 0.0
        f5 = 0.0
        f6 = 0.0
        f7 = 0.0

        if self.snake.x == self.fruit.x:
            if self.snake.y >= self.fruit.y:
```

```
343                    f2 = 1.0
344            if self.snake.y <= self.fruit.y:
345                f6 = 1.0
346        if self.snake.y == self.fruit.y:
347            if self.snake.x <= self.fruit.x:
348                f0 = 1.0
349            if self.snake.x >= self.fruit.x:
350                f4 = 1.0
351        for i in range(-self.gridlength, self.gridlength):
352            if ((self.snake.x+i) == self.fruit.x and
353                        (self.snake.y+i) == self.fruit.y):
354                if i >= 0:
355                    f7 = 1.0
356                if i <= 0:
357                    f3 = 1.0
358            if ((self.snake.x+i) == self.fruit.x and
359                        (self.snake.y-i) == self.fruit.y):
360                if i >= 0:
361                    f1 = 1.0
362                if i <= 0:
363                    f5 = 1.0
364
365        return torch.tensor([b0, b1, b2, b3,
366                             s0, s1, s2, s3,
367                             f0, f1, f2, f3, f4, f5, f6, f7])
368
369    def get_state4(self):
370        #capture the whole grid as 2d vector
371        #snake = -0.5
372        #fruit = 1
373        #body = -1
374        state = np.zeros((self.gridlength, self.gridlength))
375
376        #to omit any erros as the snake-head tried to leave the grid
377        if not self.snake.x >= self.gridlength or self.snake.x < 0:
378            if not self.snake.y >= self.gridlength or self.snake.y < 0:
379                state[self.snake.y][self.snake.x] = -0.5
380        state[self.fruit.y][self.fruit.x] = 1
381        for b in self.body.list:
382            state[b[1]][b[0]] = -1
383
384        state = torch.tensor(state)
385        state = torch.flatten(state)
```

```
386          state = state.float()
387          return state
388
389      def get_state5(self):
390          #captures a 5by5 2d array around the snakes head,
391          #body = -1
392          #border = -1
393          #fruit = 1
394          state = np.zeros((5,5))
395          for i in range(-2,3):
396              for j in range(-2,3):
397                  x = self.snake.x + i
398                  y = self.snake.y + j
399
400                  if (x >= self.gridlength or x < 0 or
401                      y >= self.gridlength or y < 0):
402                      state[j+2][i+2] = -1
403                  for b in self.body.list:
404                      if x == b[0] and y == b[1]:
405                          state[j+2][i+2] = -1
406                  if x == self.fruit.x and y == self.fruit.y:
407                      state[j+2][i+2] = 1
408
409          state = torch.tensor(state).float()
410          return torch.flatten(state)
411
412      def get_state6(self):
413          #same as get_state5
414          #+ the vector from the snakes head to the fruit
415          state = np.zeros((5,5))
416          for i in range(-2,3):
417              for j in range(-2,3):
418                  x = self.snake.x + i
419                  y = self.snake.y + j
420
421                  if (x >= self.gridlength or x < 0 or
422                      y >= self.gridlength or y < 0):
423                      state[j+2][i+2] = -1
424                  for b in self.body.list:
425                      if x == b[0] and y == b[1]:
426                          state[j+2][i+2] = -1
427                  if x == self.fruit.x and y == self.fruit.y:
428                      state[j+2][i+2] = 1
```

```
429
430         state = torch.tensor(state).float()
431         state = torch.flatten(state)
432
433         head_to_fruit = torch.tensor([self.fruit.x-self.snake.x,
434                                       self.fruit.y-self.snake.y]).float()
435
436         return torch.cat((state, head_to_fruit))
437
438     def get_state7(self):
439         #captures a 9by9 2d array around the snakes head,
440         #body = -1
441         #border = -1
442         #fruit = 1
443         state = np.zeros((9,9))
444         for i in range(-4,5):
445             for j in range(-4,5):
446                 x = self.snake.x + i
447                 y = self.snake.y + j
448
449                 if (x >= self.gridlength or x < 0 or
450                     y >= self.gridlength or y < 0):
451                     state[j+4][i+4] = -1
452                 for b in self.body.list:
453                     if x == b[0] and y == b[1]:
454                         state[j+4][i+4] = -1
455                 if x == self.fruit.x and y == self.fruit.y:
456                     state[j+4][i+4] = 1
457
458         state = torch.tensor(state).float()
459         state = torch.flatten(state)
460
461         head_to_fruit = torch.tensor([self.fruit.x-self.snake.x,
462                                       self.fruit.y-self.snake.y]).float()
463
464         return torch.cat((state, head_to_fruit))
465
466     def get_state8(self):
467         #same as get_state7
468         #+ the vector from the snakes head to the fruit
469         state = np.zeros((9,9))
470         for i in range(-4,5):
471             for j in range(-4,5):
```

```
472                     x = self.snake.x + i
473                     y = self.snake.y + j
474
475                     if (x >= self.gridlength or x < 0 or
476                         y >= self.gridlength or y < 0):
477                         state[j+4][i+4] = -1
478                     for b in self.body.list:
479                         if x == b[0] and y == b[1]:
480                             state[j+4][i+4] = -1
481                     if x == self.fruit.x and y == self.fruit.y:
482                         state[j+4][i+4] = 1
483
484         state = torch.tensor(state).float()
485         state = torch.flatten(state)
486
487         head_to_fruit = torch.tensor([self.fruit.x-self.snake.x,
488                                       self.fruit.y-self.snake.y]).float()
489
490         return torch.cat((state, head_to_fruit))
491
492     def env_step(self, action):
493         #takes in an action and puts it through the environment
494         self.snake.act(action)
495         self.body.update(self.snake.changelist, self.snake)
496         if self.collisionfruit():
497             self.body.grow(self.snake)
498             self.body.update(self.snake.changelist, self.snake)
499             self.score += 1
500             if self.score == (self.gridlength**2 - 1):
501                 #incase the agent ever manages to fill the whole grid
502                 return 100
503             self.check_fruit_spawn()
504             return self.r_fruit
505         elif self.collisionborder():
506             return self.r_collision
507         elif self.collisionself():
508             return self.r_collision
509         else:
510             return self.r_step
511
512     def env_reset(self):
513         #resets the environment and return fresh state
514         self.score = 0
```

```
515         self.snake = Snake(self.snakex, self.snakey)
516         self.fruit = Fruit(self.gridlength,
517                             spawnset=self.fruitspawnset,
518                             spawnx=self.fruitx, spawny=self.fruity)
519         self.body = Body(0)
520         return self.get_state()
521
522     def env_store_score(self):
523         #used to save a list of scores when the agent is training
524         self.scorelist.append(self.score)
525         self.score = 0
526
527
528     def render(self, size, textcolor, record=False, path=None):
529
530         #can be used to render the environment
531         import pygame
532         pygame.init()
533         pygame.font.init()
534         self.screen = pygame.display.set_mode(((self.gridlength*size),
535                                                 (self.gridlength*size)))
536         pygame.display.init()
537         self.screen.fill(self.background)
538         for b in self.body.list:
539             x = b[0]*size
540             y = b[1]*size
541             pygame.draw.rect(self.screen, self.bodycolor,
542                                             ((x,y),(size,size)))
543         pygame.draw.rect(self.screen, self.snakecolor,
544                     (self.snake.x*size, self.snake.y*size, size, size))
545         pygame.draw.rect(self.screen, self.fruitcolor,
546                     (self.fruit.x*size, self.fruit.y*size, size, size))
547         pygame.display.update()
548         if record:
549             pygame.image.save(self.screen,path)
550
551     def save_svg(self,path):
552         #can be used to save an svg of the current game state
553         d = draw.Drawing(self.gridlength*10,
554                         self.gridlength*10, origin=(0,0))
555         b = draw.Rectangle(0,0,self.gridlength*10,
556                             self.gridlength*10, fill='#485B57')
557         s = draw.Rectangle(self.snake.x*10,self.snake.y*10,10,10,
```

```
558                          fill='#197416',stroke='black', stroke_width=0.02)
559         f = draw.Rectangle(self.fruit.x*10,self.fruit.y*10,10,10,
560                          fill='#E92323',stroke='black', stroke_width=0.02)
561         d.append(b)
562         d.append(s)
563         d.append(f)
564         for b in self.body.list:
565             r = draw.Rectangle(b[0]*10,b[1]*10,10,10,
566                          fill='#32B145',stroke='black', stroke_width=0.02)
567             d.append(r)
568         d.saveSvg(path)
569
570     def renderreplay(self,gridlength, size,
571                     snakex, snakey, fruitx, fruity, bodylist):
572         #can be used to watch an episode of the agent's training history
573         import pygame
574         screen = pygame.display.set_mode(((gridlength*size),
575                                           (gridlength*size)))
576         pygame.display.init()
577         screen.fill((96,96,96))
578         for b in bodylist:
579             x = b[0]*size
580             y = b[1]*size
581             pygame.draw.rect(screen, self.bodycolor, ((x,y),(size,size)))
582         pygame.draw.rect(screen, self.snakecolor,
583                                 (snakex*size, snakey*size, size, size))
584         pygame.draw.rect(screen, self.fruitcolor,
585                                 (fruitx*size, fruity*size, size, size))
586         pygame.display.update()
```

Code 8: Source Code for the Snake Environment

# I  Source Code: Snake Main Script

```python
#############################################################
# This code is a demonstration how to use the SnakeTrainer
# To execute, please place snake_agent.py, snake_env.py,
# snake_trainer.py and snake_main.py in the same directory
#############################################################


import snake_trainer
import numpy as np

#initializing the trainer
trainer = snake_trainer.SnakeTrainer(agent_id=0, updates_per_step=32,
                    gamma=0.98, lr=0.00005, epsilon=1, epsilon_decay=0.9998,
                    epsilon_min=0.01, statetype=8, nr_hiddenlayers=1,
                    nr_neurons=[128], max_mem_length=10000, gridlength=9,
                    r_fruit=50, r_collision=-50, r_step=0, lr_decay=False,
                    lr_milestones=[80000],
                    lr_gamma=0.5, render_bool=False)


#training the agent
trainer.train_n_episodes(n=200,gridlength=3,epsilon_reset=False)
trainer.train_n_episodes(n=400,gridlength=5,epsilon_reset=False)
trainer.train_n_episodes(n=800,gridlength=7,epsilon_reset=False)
trainer.train_n_episodes(n=1600,gridlength=9,epsilon_reset=False)

#watching a replay of his training games
trainer.watch_replay(episodes=[150,2100],gridlengths=[3,9])

#testing the agent
scores = trainer.test_n_episodes(n=10,gridlength=9,render=False,record=False)
print(f'average: {np.average(scores)}')
print(f'high score: {max(scores)}')
```

Code 9: Source Code for the Snake Main Script