



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

BACHELOR THESIS

Una Adilović

**Playing a 3D Tunnel Game Using
Reinforcement Learning**

Department of Software and Computer Science Education (KSVI)

Supervisor of the bachelor thesis: Adam Dingle, M.Sc.

Study programme: Computer Science

Study branch: Artificial Intelligence

Prague 2023

I declare that I carried out this bachelor thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date
Author's signature

I would like to express my heartfelt appreciation to my supervisor, Adam Dingle, for his invaluable guidance and support throughout my academic journey, beginning with our first Programming 1 class and continuing through to the completion of this project. His patience and assistance were greatly appreciated and played a significant role in the success of this work.

In addition, I am deeply thankful to my mother for her unwavering belief in me and to my grandparents who made it possible for me to be here. Their contributions are greatly appreciated and will never be forgotten.

Title: Playing a 3D Tunnel Game Using Reinforcement Learning

Author: Una Adilović

Department: Department of Software and Computer Science Education (KSVI)

Supervisor: Adam Dingle, M.Sc., Department of Software and Computer Science Education (KSVI)

Abstract: Tunnel games are a type of 3D video game in which the player moves through a tunnel and tries to avoid obstacles by rotating around the axis of the tunnel. These games often involve fast-paced gameplay and require quick reflexes and spatial awareness to navigate through the tunnel successfully. The aim of this thesis is to explore the representation of a tunnel game in a discrete manner and to compare various reinforcement learning algorithms in this context. The objective is to evaluate the performance of these algorithms in a game setting and identify potential strengths and limitations. The results of this study may offer insights on the application of discrete tabular methods in the development of AI agents for other continuous games.

Keywords: tunnel game, reinforcement learning, artificial intelligence, algorithms

Contents

Introduction	3
1 Game Design	4
1.1 Player and Movement	4
1.2 Obstacles	5
1.2.1 Traps	5
1.2.2 Bugs	5
1.2.3 Viruses	6
1.3 Additional Features	6
1.4 Score Count and Winning	7
2 Implementation of the Game	8
2.1 The top-level organization	8
2.2 Game	9
2.3 Hans	9
2.4 Tunnels	10
3 Structure of the Experimental Setting	13
3.1 State	13
3.2 Command-line options	14
3.2.1 Command-line option descriptions	15
3.2.2 Running the program	16
3.3 Main	18
4 Applied Algorithms	20
4.1 Monte Carlo	21
4.2 Temporal Difference Learning	22
4.2.1 SARSA	23
4.2.2 Q-learning	23
4.2.3 Expected SARSA	24
4.2.4 Double Q-learning	24
5 Implementation of the Agents	26
5.1 Hierarchy	26
5.2 Simple Agents	27
5.3 Learning Agent	27
5.3.1 Common parameters and behaviours	27
5.3.2 Monte Carlo Agent	28
5.3.3 TD Agent	29
6 Experiments	32
6.1 Conducting the experiments	32
6.2 Plotting	33
6.3 Interesting behaviours	35
6.4 Individual traps environment	37
6.5 Traps environment	41

6.6	Tokens environment	41
6.7	Bugs environment	41
6.8	Viruses environment	41
6.9	Full game environment	41
Conclusion		42
Bibliography		43
List of Figures		44
List of Tables		45
List of Abbreviations		46
A	Attachments	47
A.1	First Attachment	47

Introduction

Reinforcement learning is a subfield of machine learning that aims to train agents to make decisions that will maximize a reward signal [Sutton and Barto, 2018]. This approach has been widely applied in the field of artificial intelligence, particularly in the context of training agents to play games. In a game setting, an agent’s actions can be evaluated based on their impact on the agent’s score or likelihood of winning. Through the process of reinforcement learning, the agent learns to make strategic decisions that maximize its reward by receiving positive reinforcement for good moves and negative reinforcement for suboptimal moves. This allows the agent to adapt and improve its performance over time as it plays the game. Research on reinforcement learning in games has demonstrated its effectiveness in a variety of contexts, including board games, video games, and real-time strategy games.

In the field of artificial intelligence, games can be classified as either continuous or discrete based on the nature of the action space and state space. Continuous games have a continuous action space, meaning that the possible actions an agent can take are not limited to a fixed set of options, but can vary continuously within a certain range. In contrast, discrete games have a discrete action space, meaning that the possible actions are limited to a fixed set of options. Continuous games are often characterized by a high-dimensional state space, as they may involve a large number of variables that describe the game state. Discrete games, on the other hand, typically have a lower-dimensional state space, as the number of possible states is limited by the discrete action space. In general, continuous games are more challenging to model and solve than discrete games, as they require more complex decision-making algorithms and may require more computational resources.

In this thesis, we will investigate the application of reinforcement learning to train agents to play a continuous 3D tunnel game, which I designed and implemented myself for this thesis work. The continuous game environment will be discretized into a set of states, and different reinforcement learning algorithms will be applied to train agents to play the game. The goal of this study is to determine whether it is possible for any of the agents to win the whole game, and to compare the performance of different agents that use different reinforcement learning algorithms.

The results of this study will contribute to the understanding of the potential of reinforcement learning for training agents to use discrete algorithms in a naturally continuous environment, and to provide insight into the strengths and weaknesses of different reinforcement learning algorithms in this context.

1. Game Design

For this thesis work I designed and implemented a game called “Space-Run”, which involves attempting to accumulate the highest score possible by navigating through three distinct tunnels while avoiding various obstacles. The game is endless in nature, as the speed increases each time the player successfully completes all three tunnels.

It is worth noting that, in designing this game, I was inspired by the pre-existing “Tunnel Rush” (tun [2023]). “Tunnel Rush” and “Space Run” are both 3D tunnel games that involve advancing through a tunnel to avoid traps. However, there are several key differences between the two games. “Tunnel Rush” is a web-based game played in first-person perspective, while “Space Run” is a desktop-based game played in third-person perspective. “Tunnel Rush” has levels, some of which are inverted with the traps on top of the tunnel and the player outside of it. “Space Run”, since inspired by “Tunnel Rush”, also has levels, but they are all inside the tunnel and features not only traps but also creatures that the player must avoid or shoot. Additionally, “Space Run” has a computer-themed setting, with elements such as battery, bugs, and viruses.

1.1 Player and Movement

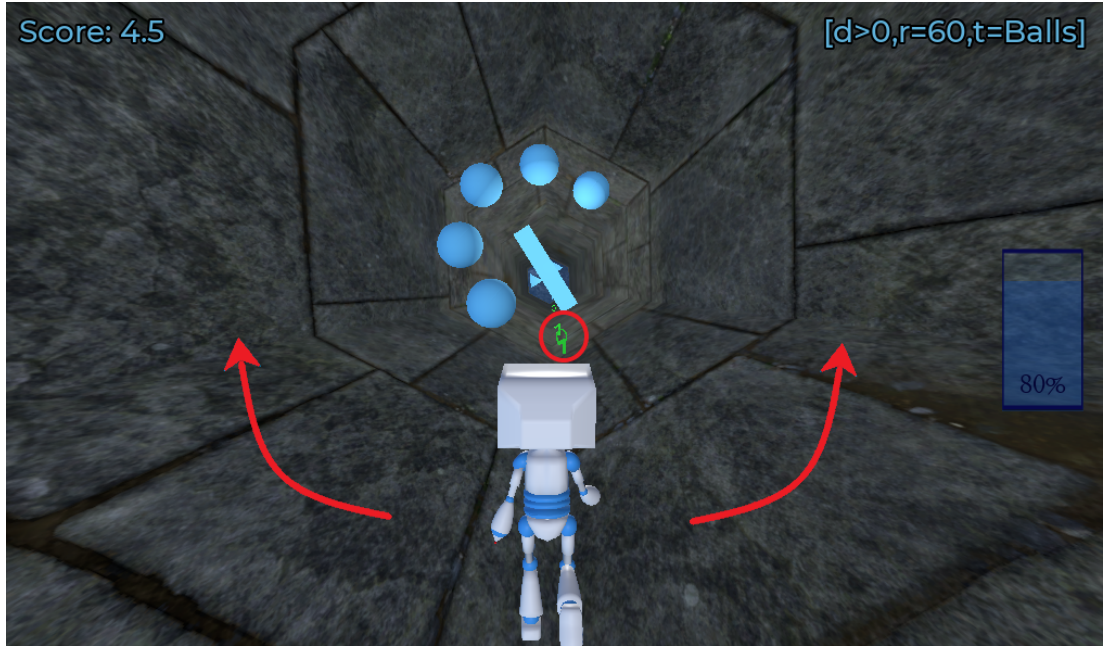


Figure 1.1: Movement

In “Space-run” the player assumes control of a character named Hans (see Figure 1.2) who continually advances at a constant speed through the tunnels. To navigate through the game, the player must use the left and right arrow keys to rotate the current tunnel and avoid obstacles (as shown in Figure 1.1).

In addition to these lateral movements, Hans also has the ability to shoot bullets (also shown in Figure 1.1) by pressing the space key, which can be used to defeat certain in-game creatures and earn a higher score. The player must utilize these abilities in order to progress through the game and achieve a high score.



Figure 1.2: Hans

1.2 Obstacles

As previously stated, the player must navigate through various obstacles in the game. These obstacles can be divided into three distinct categories, and each tunnel contains a unique subset of them. In the subsequent sections, we will delve deeper into these categories in order to better understand the challenges faced by the player.

1.2.1 Traps

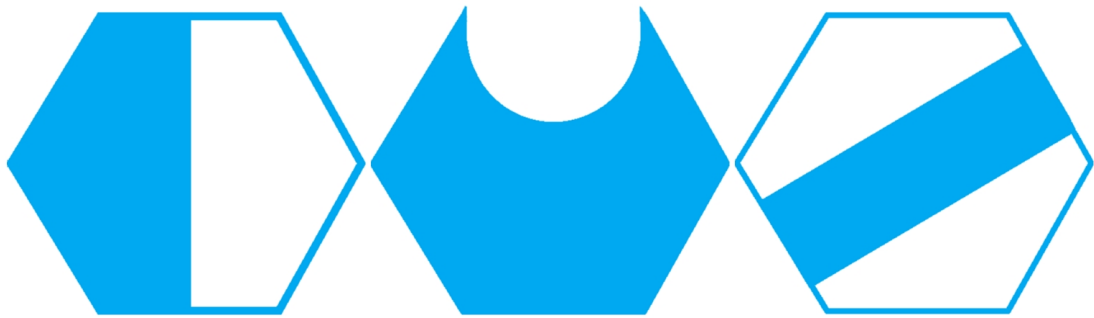


Figure 1.3: Trap examples

As depicted in Figure 1.3, a selection of the various trap types that the character Hans must avoid is presented. These traps, of which there are a total of 10, vary in their level of difficulty and can be either static or animated. These traps can be encountered in any of the three tunnels, and if the player fails to successfully evade them, they result in an instant death.

1.2.2 Bugs



Figure 1.4: Bugs

In addition to traps, the game also features bugs as an obstacle (as shown in Figure 1.4). These bugs typically appear in the second tunnel and are designed to rotate around the tunnel toward the player's position, making them more challenging to evade. However, they can still be avoided by the player. If the player chooses to engage with the bugs, they can be defeated by shooting three bullets at them. If the player collides with a bug, Hans will lose 25% of his battery life (for more information on battery life, see Section 1.3).

1.2.3 Viruses

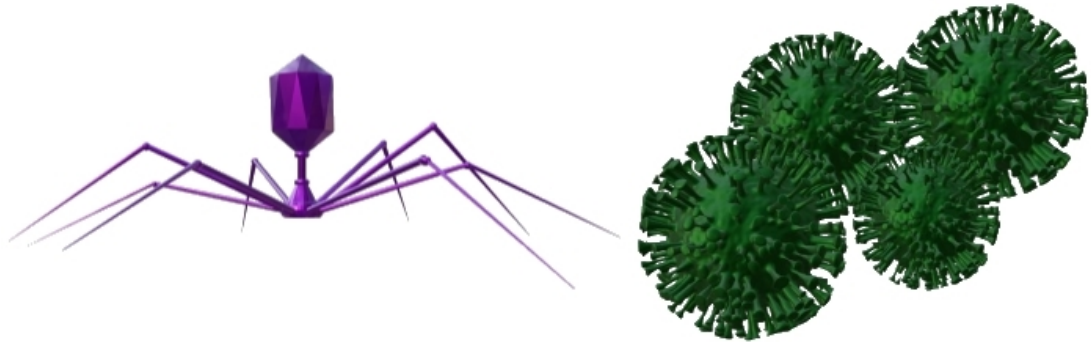


Figure 1.5: Bacteriophage and Rotavirus

The third and final type of obstacle in the game are viruses (illustrated in Figure 1.5). These viruses are typically found in the third tunnel and, similar to bugs, can be eliminated through the use of three bullets. They also, just like bugs, rotate around the tunnel toward the player's position. Bacteriophage, a subtype of virus, will result in an instant death if the player comes into contact with them. Rotaviruses, on the other hand, will cause the player's character to become sick for a brief period of time. During this illness, it is crucial for the player to avoid coming into contact with another Rotavirus, as this will result in the end of the game.

1.3 Additional Features

There are several other features of the game that are worth mentioning. One of the most significant of these is the battery life of the player's character, Hans, which is displayed on the right side of the screen (as shown in Figure 1.6). As Hans is designed to resemble a computer, it is necessary for him to recharge his battery throughout the game by collecting energy tokens (Figure 1.6). This will fully restore his battery capacity. There are three main ways in which Hans can lose battery life: running causes a constant reduction of 1% every 0.5 seconds, each bullet shot costs 1% of the battery life, and coming into contact with a bug results in a reduction of 25% (as described in Section 1.2.2). If the battery reaches 0%, Hans will die and the game will end.

Finally, it should be noted that upon successfully navigating through all three types of tunnels, the game will increase in speed and the player will once again encounter the same tunnels, looping through them indefinitely until the player loses.



Figure 1.6: Battery and Energy Token

1.4 Score Count and Winning

The score of the game is based on the length of time that the player is able to survive. Additionally, each time a player successfully shoots down a bug or virus, their score increases by 10 points. As previously mentioned, the game is designed to be played indefinitely, but for the purpose of this study, we have set the game to be considered won after an agent successfully completes nine tunnels, reaching level 10.

2. Implementation of the Game

“Space-run” was developed using the Godot Engine (version v3.2.3.stable), an open-source game engine licensed under the MIT License. It is a cross-platform tool that offers a range of features for game development, including a visual scripting language, 2D and 3D graphics support, and a powerful physics engine. The Godot Engine utilizes a node-based architecture, where nodes are organized within scenes that can be reused, instanced, inherited, and nested. This structure allows for efficient project management and development within the engine. The game was written entirely in GDScript, the primary scripting language of the Godot Engine (Linietsky [2021]).

In addition to using the Godot Engine, I also utilized Blender (version 6.2.0) (Roosendaal [2023]) for creating and animating the characters in the game. Blender is popular open-source 3D modeling and animation software that offers a range of features for creating detailed and realistic characters. The characters were then imported into the Godot Engine using the .glTF 2.0 (Khronos [2023]) file format, which is a widely supported file format for exchanging 3D graphics data.

The game can be run on any platform provided within the Godot engine, and its source code, and the source code for the whole thesis can be found online (Adilović [2023]).

2.1 The top-level organization

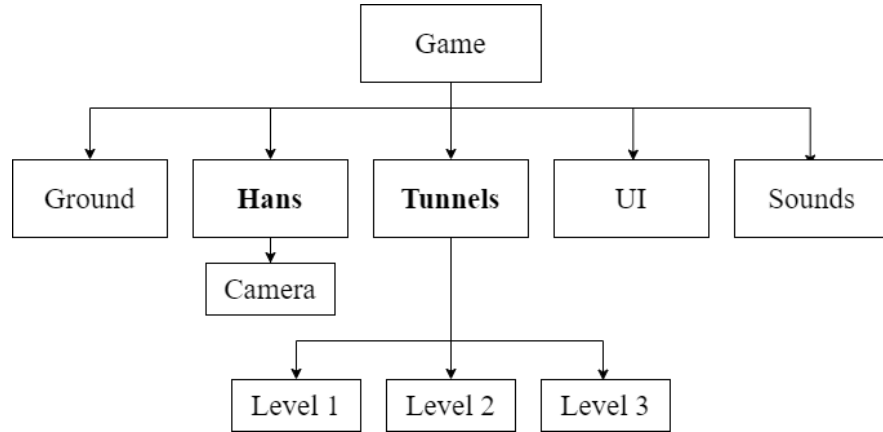


Figure 2.1: Structure of Game.tscn

The main scene for the game, referred to as `Game.tscn`, is depicted in Figure 2.1. It includes several nodes, including Ground, UI, Sounds, Game, Hans, and Tunnels. The Ground node is a `CSGBox`¹ that serves as the ground in the game, while the UI and Sounds nodes handle the user interface and audio aspects, respectively. The Game, Hans, and Tunnels nodes contain the majority of the game’s functionality. Specifically, the Game node manages the overall gameplay,

¹A `CSGBox` is a 3D object that represents a box with a Constructive Solid Geometry (CSG) shape.

the Hans node controls the player character, and the Tunnels node manages the movement and appearance of the tunnels.

For a more in-depth understanding, let us examine some of the core aspects of the game in the following sections.

2.2 Game

The script for the Game node is the initial point of the game session and includes both the `_start()` and `_game_over()` methods. It also serves as a link between the game and the agent environment described in Chapter 3, and as such includes all of the necessary set methods for the agent environment. These methods allow for communication between the game and the agent environment, enabling the agent to interact with and influence the game.

The following text describes the core functionalities of the main methods within `Game.gd`:

- The `_ready()` function is called at the start of the game’s execution and, after setting up the environment, it triggers the `_start()` function. This function initiates the gameplay and sets the necessary conditions for the game to proceed.
- As described in more detail in Chapter 3, the user can specify environment parameters and a starting level for the agent through the command line. These parameters determine the obstacles that the player will face and the starting position of the player character, Hans. The `_start()` function incorporates these parameters into the obstacle arrays and positions Hans accordingly. The function also generates the obstacles for the designated starting level. The creation and deletion of obstacles during gameplay is discussed in Section 2.4 of this chapter.
- The `_game_over` function manages the end of the game and sends a signal to the top-level script, `Main.gd` (described in Chapter 3), indicating that the game has ended. It also provides `Main.gd` with the necessary information about the game’s status and outcome.

2.3 Hans

The next node we want to examine is Hans. While `Hans.tscn` is a scene with the main character and its necessary animations, what interests us more is the `Hans.gd` and its key components.

```
func physics_Process(delta):
    tunnels.deleteObstacleUntilX(...)

    if translation.x < new_trap:
        createNewTrap()
```

```

score._on_Meter_Passed() # update score

velocity = Vector3.LEFT * speed
velocity = moveAndSlide(velocity)

checkCollisions()

tunnels.bugVirusMovement(delta, curr_tunnel)

if isShootingButtonPressed:
    shoot()

state.updateState(...)

```

The primary function within `Hans.gd` is the `_physics_process()`, which is called on every tick of the game. It handles the main aspects of the player character through the use of various methods and functions. These include deleting passed obstacles, creating new obstacles every 50 meters, updating the score, handling the movement of the player character, bugs and viruses, and determining the current state of the player. The state label, which is displayed on the upper right corner of the screen (as shown in Figure 2.2), is the primary information that agents receive when making decisions about their next move, as described in Chapter 3. The `_physics_process()` function also handles collisions and shooting if the player chooses to do so. Overall, this function plays a crucial role in the gameplay and management of the player character.

It is also worth noting that this script handles the movement of the tunnels to the back as Hans passes them, with the first tunnel being moved to be after the third one. This feature allows for the game to be infinite, as the tunnels are constantly cycled and reused.

2.4 Tunnels

The Tunnels node, which is a child of the main scene in the game tree, contains three child nodes of the Spatial type² (level1, level2, and level3) and each of these nodes includes a CSGTorus³ node, which represents the physical appearance of the tunnels. Obstacles are added to the appropriate level node as instances. The `Tunnels.gd` script, which is attached to the Tunnels node, handles many of the previously mentioned functions such as obstacle creation and deletion and tunnel rotation. In the following code snippets, we will examine the `Tunnels.gd` script in greater detail.

²A Spatial node is a type of node that represents a 3D object or transformation in the game world. It is a versatile node that can be used to create and manipulate 3D objects, including meshes, materials, and lighting. Spatial nodes are often used as the root node for 3D objects in a scene, and they can be nested inside other Spatial nodes to create hierarchical transformations.

³A CSGTorus node is a type of 3D object that represents a torus shape in the game world.



Figure 2.2: State

```
func physics_process(delta):
    var move = game.agent.move(...)
    if move[0] == 1:
        tunnel = get_child(hans.get_current_tunnel())
        tunnel.rotate_object_local(Vector3.RIGHT, -
            ROTATE_SPEED * delta)
    elif move[0] == -1:
        tunnel = get_child(hans.get_current_tunnel())
        tunnel.rotate_object_local(Vector3.LEFT, -
            ROTATE_SPEED * delta)

    if hans != null: # If it is not instanced we can't
        call the function
        hans.switch_animation(move[1] == 1) # Shoot if
            necessary
end
```

The `_physics_process()` function within the `Tunnels.gd` script serves as the primary connection between the agent and the game. As shown in the provided code, the function retrieves the next move from the agent and rotates the tunnel accordingly, potentially including shooting as well. It should be noted that, in the case of the Keyboard agent, function `move` returns users input from the keyboard.

```
func create_first_level_traps(tunnel):
    var level = tunnel
    var num_of_traps = rand_range(
        MIN_TRAPS_PER_TUBE, MAX_TRAPS_PER_TUBE)
```

```

x = position of the first trap

for n in range(num_of_traps):
    # update x position
    x -= rand.randi_range(TRAP_RANGE_FROM,
                          TRAP_RANGE_TO)
    if x is outside the tunnel:
        break
    create_one_obstacle(level, x)

```

The function depicted in the code above serves to generate obstacles in the starting tunnel. By periodically creating traps in the tunnel ahead, the game is able to prevent lag caused by an excessive number of objects existing simultaneously. For that reason, this function is used only once, at the beginning of the game.

```

func CreateOneObstacle(level, x):
    var scene =
        pick_which_kind_of_obstacle_will_be_added
    var tunnel =
        get_the_level_we_are_making_traps_for
    var i = randomly_pick_an_obstacle

    # make an instance
    var obstacle = scene[i].instance()
    obstacle.translation.x = x

    tunnel.add_child(obstacle)
    rotate_obstacle(obstacle)

```

The tunnels are positioned along the x axis, and this function allows for the creation of obstacles within them at specific x positions and a random rotation.

```

func delete_obstacle_until_x(level, x):
    var tunnel = get_current_tunnel()
    for obstacle in tunnel.get_children():
        if obstacle is an obstacle type:
            if obstacle.translation.x > x:
                obstacle.queue_free()
        else:
            return

```

As previously mentioned, by dynamically deleting passed obstacles, the game is able to maintain a stable performance and avoid overloading the system.

3. Structure of the Experimental Setting

To train an agent on a specific environment, the user must utilize a command-line interface. There are various options available to cater to the user’s needs. This chapter will outline all of the provided options and how they are encoded. However, we will first consider how the game represents discrete states.

3.1 State

The majority of the implemented agents in the game utilize the concept of state to facilitate learning. The agent will determine its next action based on the current state in which it finds itself. By dividing the game into discrete states, we are able to utilize tabular method algorithms to train an agent to play this continuous game. Once one obstacle is passed, the value of the state indicated refers to the next one. Each state is represented by a triplet consisting of distance, rotation, and next obstacle type.

It is important to note that the game uses meters as the standard unit of distance measurement in Godot. This unit is used to represent the size, position, and movement of objects in the game world. It is worth mentioning that one meter in Godot is equal to the size of a standard cube. With that in mind, players can visualize Hans’ approximate height of 12 meters and starting speed of 35 meters per second. Additionally, each tunnel in the game has a width and height of approximately 45 meters and a depth of 2800 meters.

The distance value indicates the distance of the agent from the next obstacle in meters. For example, if we set the `dist`s parameter (indicated by the user, see 3.2.1) to 2, there are two possible distance values for the state: `>50`, `>0`, indicating that the agent is more than 50 meters away from the obstacle and somewhere between 1 and 50 meters away from the obstacle, respectively. If the `dist`s parameter is set to 1, the distance value remains constant at `>0`. In the Chapter 6 we will see that all types of environments in this game are able to learn when the `dist`s parameter is equal to 1.

The rotation parameter divides the obstacle circumference into `rots` number of intervals (see 3.2.1). As the tunnel rotates, the state label will indicate the rotation value to which Hans is aligned at the moment (in Figure 3.1 Hans is aligned with rotation value 240). If the `rots` parameter is set to 360, this would correspond to the number of degrees in a circle and result in 360 possible rotations. However, the obstacles in this game do not require such a high number of rotations and agents can be trained to avoid most obstacles using fewer than 10 rotations.

Finally, the type parameter indicates the type of obstacle that Hans must avoid. Each obstacle has its own unique string representation, which allows the agent to learn to recognize the safe rotation for different types of obstacles. Using this information, along with the distance parameter, the agent can decide whether to move left, right, forward, or shoot in combination with any of these actions. To further explain this notion let’s look back at Figure 3.1. At first glance, it may seem that a rotation of 240 would allow Hans to easily pass through the

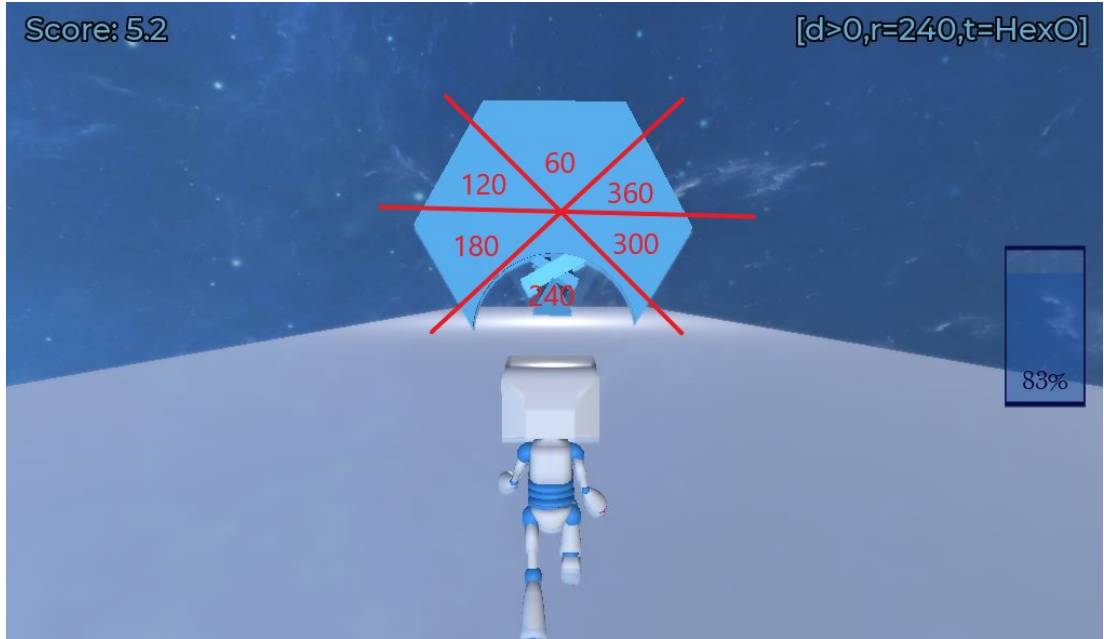


Figure 3.1: Rotations when `rots = 6`

obstacle. However, upon closer inspection, it becomes clear that only a specific portion of the 240 rotation value is safe. The agent does not know the difference between being at any point of the 240 rotation and thus can easily make a fatal mistake. For this obstacle, indeed, there would be more `rots` needed in order for Hans to have at least one safe rotation value (meaning that the whole piece of the obstacle that rotation value covers is considered safe).

In Chapter 6, a noteworthy concept is presented where the agent can recognize a safe edge between two rotations, even if it is not provided with a completely safe rotation. By quickly shifting between these rotations, the agent is able to successfully navigate through the obstacle.

3.2 Command-line options

<code>n=int</code>	number of games
<code>stoppingPoint=int</code>	stop after the agent wins this many consecutive games
<code>agent=string</code>	name of the agent
<code>level=int</code>	number of the level to start from
<code>env=[string]</code>	list of obstacles that will be chosen in the game
<code>shooting=string</code>	enable or disable shooting
<code>dists=int</code>	number of states in a 100-meter interval
<code>rots=int</code>	number of states in 360 degrees rotation
<code>agentSeedVal=int</code>	seed value for the random moves the agent takes
<code>database=string</code>	read the data for this command from an existing file and/or update the data after the command is executed
<code>ceval=bool</code>	performs continuous evaluation
<code>debug=bool</code>	display debug print statements
<code>options</code>	displays options

3.2.1 Command-line option descriptions

In this section, a more comprehensive clarification for the table presented earlier is provided. It is important to note that any of the options listed can be left out, as they all have default values assigned to them. In the event that no options are specified, a regular game with the `Keyboard` agent will be executed.

n

Number of games the agent will train on in this session. The default is 100.

stoppingPoint

If the agents manages to win this many consecutive games, the experiment is stopped. This is used to avoid long experiments if the agent has already learned the right policy and would simply keep winning. The default value is 25.

agent

Name of the desired agent.

Options: [`Keyboard`, `Static`, `Random`, `MonteCarlo`, `SARSA`, `QLearning`, `ExpectedSARSA`, `DoubleQLearning`]

Sub-options (only for the learning agents):

[`gam` (range [0,1]), `eps` (range [0,1]), `epsFinal` (range [0,1]), `initOptVal` (range [0,∞))]

Example usage: “`agent=MonteCarlo:eps=0.1,gam=0.2`”

The meaning of the suboptions is explained in Chapter 4.

level

Number of the level to start from. Default value is 1.

Options: [1, . . . , 15]

Note: after the 15th level, the agent is considered to have won the game.

env

List of obstacles that will be chosen in the game,

Options (any subset of): [`Traps`, `Bugs`, `Viruses`, `Tokens`, `I`, `O`, `MovingI`, `X`, `Walls`, `Hex`, `HexO`, `Balls`, `Triangles`, `HalfHex`, `Worm`, `LadybugFlying`, `LadybugWalking`, `Rotavirus`, `Bacteriophage`]

Note: if this parameter is not included, the environment will contain all available obstacles (i.e. the full game).

Example usage: “`env=HexO,I,Bugs`”

shooting

Enable or disable shooting.

Options: [`enabled`, `disabled`, `forced`]

This option is disabled by default or if the environment does not have any bugs or viruses.

dists

Number of states in a 100-meter interval.

This parameter is part of the state label and typical options range from 1 to 3. Default value is 1.

rots

Number of states in 360 degrees rotation. This parameter is part of the state label and the minimum viable option is 6. This is also the default value.

agentSeedVal

Seed value for the agent.

This parameter is used to produce multiple experiments with different random actions to verify if the agent can really learn under certain conditions or if it has simply been “luck” with the random moves.

database

Read or write data for this command from/to a file.

Options: [read, write, read.write]

Note: these files are typically used to start another session of the agent’s training from the last point of the previous session, to run a game with visuals and observe the agent’s performance, or for plotting the results. This option does not affect the Keyboard, Static, or Random agents. Default option for this parameter is to neither read nor write.

ceval

Performs continuous evaluation.

This parameter indicates that after each training game, a test game will be played using only the policy(s) learned thus far. For example, if the user specifies “n=100”, a total of 200 games will be executed, with 100 of them being training games and the remaining 100 being test games. This allows for the assessment of the agent’s progress and performance during the training process. Options: [true,false (default)]

debug

Display debug print statements. Options: [true,false (default)]

options

Displays all of the mentioned options.

3.2.2 Running the program

There are several possibilities for running the game from the command line. In addition to various combinations of the options listed above, the user has the choice of running an experiment with or without the graphical interface. If they

opt for the first possibility, the window will open and the game will be played at its normal speed. On the other hand, if the experiment is run without graphics, it will be over 200 times faster and the output will only be displayed in the terminal. The program achieves this speedup by hiding the CSG geometry in every node and performing a few other tricks. The computational power required to perform union, intersection, etc. on the CSG shapes is quite high and thus by not performing those calculations a game can run much faster. Of course, these shapes are not necessary for the experiments run without graphics, since the collision shapes are the ones that play a role in determining what happened in the game¹.

It is important to acknowledge that the experiments conducted in this study are entirely reproducible owing to the predetermined seed values for all the random variables. However, it should be noted that the values produced by a seed can differ across different versions of the Godot Engine. In order to obtain identical results to the experiments detailed in Chapter 6, it is recommended to use Godot Engine v3.2.3.stable.

To run the program in the command line, the user should add the directory containing the Godot executable to the PATH environment variable. This will allow them to start the application from the command line simply by entering the command `godot` while inside the same directory as the `project.godot` file.

By default, running the program in this manner will launch a normal game with the graphical interface and the `Keyboard` agent. However, the user can customize their experiment by using a combination of the options listed above. For example:

```
$godot database=write agent=SARSA:initOptVal=100.0,eps=0.3 env=Hex0 n=10 dists=1 rots=8
```

Alternatively, the user may choose to train the agent faster by disabling the graphical interface and increasing the speed of the program. This can be achieved by modifying the previous command as follows:

```
$godot --no-window --fixed-fps 1 --disable-render-loop database=write agent=SARSA:initOptVal=100.0,eps=0.3 env=Hex0 n=10 dists=1 rots=8
```

To view a list of available options, the user can simply enter the command `godot options`.

¹The collision shapes refer to the shapes that are used to define the physical bounds of an object for the purpose of collision detection.

3.3 Main

The `Main.tscn` scene is the top level scene in the game and consists of a single Node type node. The script attached to this node, `Main.gd`, is responsible for ensuring that all options specified in the command line (as discussed in Section 3.2) are executed correctly. This script is the starting point of the training and handles the initialization and execution of one or more game sessions. There are several key functions within the `Main.gd` script that are worth discussing in more detail.

```
func _ready():
    var unparsed_args = OS.get_cmdline_args()
    if unparsed_args.size() == 1 and unparsed_args
        [0] == 'options':
        display_options()

    ...      # parse args

    if set_param(args) == false:
        display_options()
    else:
        instance_agent()
        build_filename()
        if not agent_inst.init(...):
            print('Something went wrong,
                please try again')
            display_options()
        play_game()
```

The `_ready()` function is the starting point of the program when run from the command line. It is responsible for parsing all of the arguments and checking their validity. If any issues are encountered, the program will display options and terminate. If the arguments are valid, an agent will be instantiated and initialized. In cases where everything is in order, first game will be played by calling the `play_game()` function.

```
func play_game():
    if agent == 'Keyboard' and VisualServer.
        render_loop_enabled:
        ... # play a regular game
    elif n > 0:
        n -= 1
        game = game_scene.instance()
        set_param_in_game()
        agent_inst.start_game(is_eval_game)
```

```
else:
    agent_inst.save(write)
    print_and_write_ending()
```

The `play_game()` function is called each time a game is played. Firstly, it will check if the agent selected was the `Keyboard` agent, and if so, the program will start one game session where the user has control of Hans. Otherwise, one of the remaining agents will take over and play the specified number of games (defined by the “n” parameter). If n games have already been played, the program will terminate after performing the last few tasks needed to save all of the knowledge gained from this particular session. Otherwise, a single game will be executed and number of games left decreased.

```
func on_game_finished(score, ticks, win, time):
    print_and_write_score(score, win)
    agent_inst.end_game(score, time)
    play_game()
```

The `game_over()` function is called when the game emits a signal indicating that it has finished. Upon execution, this function outputs the necessary information, updates the agent through the `end_game()` function, and then calls the `play_game()` function to continue the game session.

4. Applied Algorithms

In this chapter, we will discuss the algorithms employed to create agents for the game, which broadly fall into the categories of Monte Carlo methods and Temporal Difference (TD) Learning. In essence, the agents aim to maximize their reward by selecting actions that yield the greatest possible benefit. To comprehend the functioning of these algorithms, it is necessary to introduce several fundamental concepts.

A **state** represents the current status of the environment, while an **action** denotes a decision made by the agent in response to the current state. A **reward** is a scalar value that reflects the immediate feedback received by the agent for its action in a given state. The **discount factor**, usually denoted as γ , is a value between 0 and 1 that represents how much the agent values future rewards compared to immediate rewards. An **episode** refers to a sequence of states, actions, and rewards that begins with an initial state and ends when a terminal state is reached. It represents one run or iteration of the agent interacting with the environment. The length of an episode can vary depending on the problem and the algorithm being used (e.g. in a game, an episode may correspond to a single game). A **policy** is a mapping from states to actions that determines the actions an agent takes in each state. An **ϵ -greedy policy** is a policy in which the agent selects the action that maximizes the expected reward with a probability of $(1 - \epsilon)$, while taking a random action with a probability of ϵ . Two types of value functions exist, namely **state-value** functions and **action-value** functions. The former predict the expected long-term reward of being in a particular state, while the latter predict the expected long-term reward of taking a specific action in a particular state and always following the optimal policy thereafter. Action-value and state-value functions evaluate the relative effectiveness of different actions or states, serving as a measure to determine the optimal action in a particular state. The **value of a state** varies between algorithms, and it is used to select the best action to be taken in that state.

Before looking into individual algorithms, there is one more key concept to introduce: **exploration vs exploitation**. In the field of reinforcement learning, the exploration-exploitation trade-off refers to the balancing act between discovering new information or strategies and utilizing existing knowledge to maximize reward. Exploration involves trying out different actions or strategies in order to gather more information about the environment and its rewards, while exploitation involves utilizing the information gathered to maximize reward. Finding the right balance between exploration and exploitation is crucial in reinforcement learning, as excessive use of either can result in suboptimal results. To balance out these two concepts within these algorithms two methods are used: the formerly mentioned **ϵ -greedy policy** as well as the **optimistic initial values** which is a technique which sets the initial value of all state action pairs to a high number, encouraging the agent to visit as many states as possible in order to learn their true value [Sutton and Barto, 2018].

4.1 Monte Carlo

Monte Carlo (MC) methods are a type of reinforcement learning algorithm that estimate the value of a state or action by averaging the total reward received from sample episodes. Unlike some other methods, such as dynamic programming, MC methods do not require knowledge of the transition probabilities between states or the reward function. Instead, they learn from experience by directly observing the outcomes of sample episodes.

During each episode, the agent follows its policy to select actions, receives rewards from the environment, and transitions to the next state. Once an episode terminates, the total reward received from that episode is recorded. This total reward is used to update the value estimates for each state and action that were encountered during the episode.

There are two types of MC learning: on-policy and off-policy. On-policy learning means that the agent is using the same behavior policy to collect samples as it is using to improve the value function. Off-policy learning, on the other hand, means that the agent is using a different policy to collect samples than the policy it is using to improve the value function.

One variant of on-policy MC learning is first-visit Monte Carlo. This method only considers the first time a state is visited in an episode, as opposed to all visits. The goal of using this method is to reduce variance in the value estimates and improve learning efficiency [Sutton and Barto, 2018].

To ensure exploration during learning, the ϵ -greedy policy is often used in conjunction with MC methods. Additionally, setting an initial optimistic value can encourage the agent to visit more states to learn their true values. While Monte Carlo methods are guaranteed to converge to an optimal policy with an infinite number of samples, convergence can be slow and estimates can be noisy (i.e., have high variance) with a small number of samples.

The pseudocode for the first-visit Monte Carlo can be seen in Algorithm 1 [Sutton and Barto, 2018]¹.

To clarify this and future algorithms, here are further explanations to some of the elements that might be encountered:

- \mathcal{S} - The set of all possible states.
- $\mathcal{A}(s)$ - The set of actions possible in state s .
- S_t or A_t - Specific state or action taken at time step t .
- R_t - The reward received by the agent at time step t .
- G - The actual total return that the agent received from a single episode. In other words, it is the sum of all the rewards that the agent received from the start state until the end of the episode.
- $Returns(s, a)$ - A list that stores the observed returns (i.e., sum of rewards) that are obtained from following the policy and taking action a in state s . These returns are later used to update the action-value function $Q(s, a)$ for the state-action pair (s, a) . The list is maintained for each state-action

¹Note that all of the pseudocode in this chapter is derived from the Sutton/Barto text.

pair to keep track of the returns obtained from that state-action pair across different episodes.

- $Q(s, a)$ - The expected cumulative reward an agent would receive if it takes action a while in state s and follows a certain policy thereafter. It is a function that maps a state-action pair to a scalar value. The value of $Q(s, a)$ is updated iteratively as the agent interacts with the environment and learns from experience.
- π - particular policy the agent is following. Furthermore, $\pi(S_t, a)$ represents the probability of taking action a at state S_t under a given policy π .

Algorithm 1 On-policy first-visit Monte Carlo

Require: small $\epsilon > 0$

```

1: Initialize:
2:    $Q(s, a) \leftarrow$  initial optimistic value  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
3:    $Returns(s, a) \leftarrow$  empty list,  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
4: repeat
5:   Generate an episode following policy  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
6:    $G \leftarrow 0$ 
7:   for  $t = T - 1, T - 2, \dots, 0$  do
8:      $G \leftarrow \gamma G + R_{t+1}$ 
9:     if  $(S_t, A_t)$  does not appear in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$  then
10:      Append  $G$  to  $Returns(S_t, A_t)$ 
11:       $Q(S_t, A_t) \leftarrow$  initial optimistic value + average( $Returns(S_t, A_t)$ )
12:      Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
13:    end if
14:  end for
15: until convergence

```

In Chapter 5 of this work, we will provide a detailed discussion of the specific Monte Carlo algorithm implementation employed.

4.2 Temporal Difference Learning

Temporal Difference (TD) learning is another type of reinforcement learning algorithm that is, similarly to Monte Carlo, model-free. The main idea behind it is to update the estimated value of a state or action based on the difference between the expected return and the actual return obtained from that state or action.

TD learning is similar to Monte Carlo methods in that it learns from experience by interacting with the environment and observing the rewards received. However, these methods update their estimates after every time step, rather than waiting for an entire episode to complete like in MC methods. This makes Temporal Difference learning more efficient in terms of the amount of data needed to learn a good estimate of the value function.

Like Monte Carlo methods, TD methods can also be on-policy or off-policy. In on-policy learning, the agent learns about the value of the policy it is currently

following, whereas in off-policy learning, the agent learns about the value of a different policy.

One important parameter in TD learning is the **step size** or **learning rate** (usually denoted by the symbol α), which determines the size of the update to the value estimates. A larger step size will result in faster learning, but may also make the learning process more unstable.

In this chapter, we shall introduce four distinct TD learning algorithms, namely SARSA, Q-Learning, Expected SARSA and Double Q-Learning [Sutton and Barto, 2018]. We shall illustrate that while SARSA and Q-Learning are prominent algorithms for control problems, Expected SARSA and Double Q-Learning are variations that cater to specific limitations of the original algorithms. The objective is to highlight both the commonalities and differences among them.

4.2.1 SARSA

SARSA stands for State-Action-Reward-State-Action. With this algorithm, the agent learns the value of a state-action pair $Q(S', a)$ by estimating the expected return over all possible actions from state S' . SARSA is an on-policy algorithm, meaning it learns the value of state-action pairs while following the same policy used to select actions. This makes it well-suited for control problems, where the goal is to find an optimal policy.

Algorithm 2 SARSA

```

1: Initialize:
2:    $Q(s, a) \leftarrow$  initial optimistic value  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
3: repeat
4:   Initialize  $S$ 
5:   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
6:   repeat
7:     Take action  $A$ , observe  $R, S'$ 
8:     Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
9:      $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
10:     $S \leftarrow S'$ 
11:     $A \leftarrow A'$ 
12:   until  $S$  is terminal
13: until convergence

```

4.2.2 Q-learning

Q-learning, unlike SARSA, is an off-policy algorithm that learns the value of a state-action pair $Q(s, a)$ by estimating the maximum expected return over all possible actions from state s . In other words, it learns the value of the best action in each state. Since Q-learning is an off-policy algorithm it learns the optimal action-value function regardless of the current policy being followed. This makes Q-learning more flexible in terms of exploration and can result in faster convergence to the optimal policy. However, Q-learning tends to overestimate the value

of actions in environments with high variance, which can lead to suboptimal policies. On the other hand, SARSA is more stable and less prone to overestimating the value of actions. Nevertheless, it can converge to suboptimal policies if the exploration is insufficient, and it can take longer to converge to the optimal policy compared to Q-learning.

Algorithm 3 Q-learning

```

1: Initialize:
2:    $Q(s, a) \leftarrow$  initial optimistic value  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
3: repeat
4:   Initialize  $S$ 
5:   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy).
6:   repeat
7:     Take action  $A$ , observe  $R, S'$ 
8:      $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
9:      $S \leftarrow S'$ 
10:  until  $S$  is terminal
11: until convergence

```

4.2.3 Expected SARSA

Expected SARSA is another off-policy TD algorithm that learns the value of a state-action pair $Q(S', a)$ by estimating the expected return over all possible actions from state S' , taking into account the probabilities of selecting each action according to the current policy. Expected SARSA can be seen as a compromise between SARSA and Q-learning, as it considers the value of both the current and the best action in each state. This algorithm considers all possible actions and their expected values, which makes it more robust to noisy or uncertain rewards.

Algorithm 4 Expected SARSA

```

1: Initialize:
2:    $Q(s, a) \leftarrow$  initial optimistic value  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
3: repeat
4:   Initialize  $S$ 
5:   repeat
6:     Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
7:     Take action  $A$ , observe  $R, S'$ 
8:      $Q(S, A) \leftarrow Q(S, A) + \alpha \cdot [R + \gamma \sum_a \pi(S', a) Q(S', a) - Q(S, A)]$ 
9:      $S \leftarrow S'$ 
10:  until  $S$  is terminal
11: until convergence

```

4.2.4 Double Q-learning

Double Q-learning is a variant of Q-learning that uses two action-value functions to estimate the value of each action. The two functions are updated indepen-

dently, and the final action-value estimate is the average of the two estimates. In Q-Learning, a single estimate of the action values is used to update the policy and make decisions. This means that when selecting an action in the next state, we always select the action with the highest estimated value (here we do not take into account using ϵ -greedy policy), even if that estimate is not accurate. This can result in overestimation of the true value of that action, particularly in situations where the policy is still exploring the environment. Double Q-Learning addresses the overestimation issue in Q-Learning which can lead to more accurate value estimates and better performance in some cases.

Algorithm 5 Double Q-learning

```

1: Initialize:
2:    $Q(s, a) \leftarrow$  initial optimistic value  $\forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
3: repeat
4:   Initialize  $S$ 
5:   repeat
6:     Choose  $A$  from  $S$  using policy derived from  $Q_1 + Q_2$  (e.g.,  $\epsilon$ -greedy)
7:     Take action  $A$ , observe  $R, S'$ 
8:     if  $\text{rand}() < 0.5$  then
9:        $A' \leftarrow \arg\max_a Q_1(S', a)$ 
10:       $Q_1(S, A) \leftarrow Q_1(S, A) + \alpha[R + \gamma Q_2(S', A') - Q_1(S, A)]$ 
11:    else
12:       $A' \leftarrow \arg\max_a Q_2(S', a)$ 
13:       $Q_2(S, A) \leftarrow Q_2(S, A) + \alpha[R + \gamma Q_1(S', A') - Q_2(S, A)]$ 
14:    end if
15:     $S \leftarrow S'$ 
16:  until  $S$  is terminal
17: until convergence

```

5. Implementation of the Agents

In this project, the reinforcement learning (RL) agents are designed such that their functionality is encapsulated within a top-level scene called `Main.tscn`. Within this scene, there is an instance of the selected agent, and the code makes use of several functions implemented by the agent in order to interact with the environment. Specifically, the required functions include: `move()`, `init()`, `start_game()`, `end_game()`, `save()`.

The purpose of most of these functions is self-explanatory. `init()` and `save()` are used to initialize and save the agent's internal state, respectively, and are called only once per experiment. `start_game()` and `end_game()` are called at the beginning and end of each episode, while `move()` is called by the `Tunnels.gd` script, and it is in this function that the agent makes a decision about which action to take based on the current state and score. The remaining functions pertain to the internal structure of the agent and are not relevant to the reader.

5.1 Hierarchy

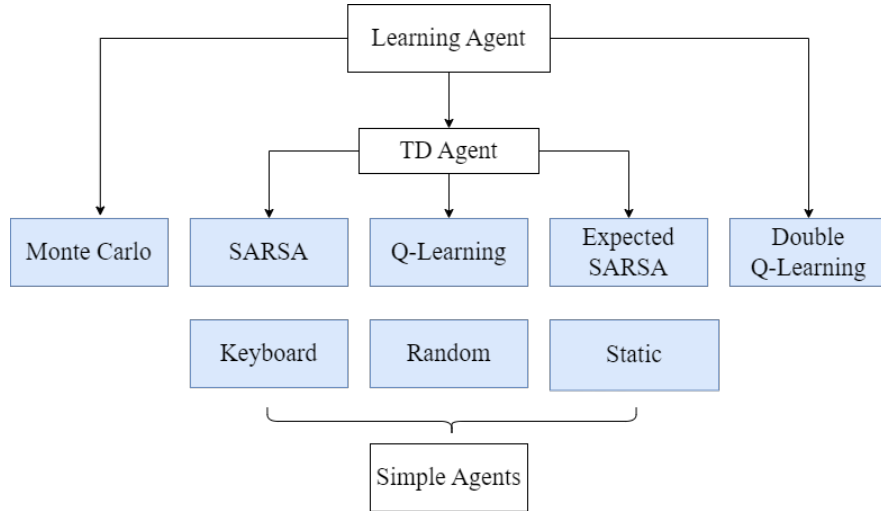


Figure 5.1: Agents hierarchy inside the project

In the current design of the game, there are a total of 8 agents implemented, 5 of which utilize some form of reinforcement learning algorithm. These RL agents share a common superclass called `LearningAgent`, while 3 of them are further subclassed under the `TDAgent` class (see Figure 5.1). As previously discussed, the RL algorithms can be broadly divided into two categories: Monte Carlo methods and temporal difference (TD) learning. The TD algorithms differ only in their update function, and so it was deemed appropriate to group them under the same superclass. However, the Double Q-Learning, which utilizes separate policies and requires additional modifications, was implemented as a separate subclass of the `LearningAgent`. The implementation details of these agents will be further elaborated upon in the subsequent sections.

Before proceeding further, it is crucial to emphasize that, as previously mentioned when describing the state in Chapter 3, the agent has the ability to per-

form movements in various directions, including moving right, forward, or left, and each of these movements can be performed in combination with shooting. Therefore, on each time step, the move function of the agent returns a list of two elements: the first element signifies the movement direction, where a value of -1 represents left, 0 represents forward, and 1 represents right; while the second element determines whether the agent will shoot or not, with a value of 1 indicating yes and 0 indicating no for shooting.

5.2 Simple Agents

To facilitate testing of the game environment, several simple agents were implemented. These agents serve as baseline models and are used to ensure that the environment is functioning as intended before more sophisticated RL agents are developed. There are three simple agents in total: a “Keyboard agent” that receives input from the player via the keyboard, a “Static agent” that always chooses the forward action without shooting, and a “Random agent” that chooses a random action at each time step.

5.3 Learning Agent

This class serves as a base class for all the reinforcement learning agents in this project. It provides a set of shared functions and features that are used by all agents, such as reading and writing data to a file and debugging statements. In terms of decision-making, these agents all follow an ϵ -greedy policy, whereby they select the action with the highest value for a given state with a certain probability, or randomly choose any action with the remaining probability. Each of the subclasses of the **Learning Agent** class then implements specific code that is unique to that particular agent.

Algorithm 6 Choosing an action using ϵ -greedy policy

```

1: function CHOOSE_ACTION(action)
2:   epsilon_action  $\leftarrow$  false
3:   if not is_eval_game then ▷ Used for continuous evaluation
4:     epsilon_action  $\leftarrow$  rand.randf_range(0, 1) < EPSILON
5:     if epsilon_action then
6:       action  $\leftarrow$  ACTIONS[rand.randi_range(0, len(ACTIONS) - 1)]
7:     end if
8:   end if
9:   return action
10: end function

```

5.3.1 Common parameters and behaviours

In regards to the present implementation of the reinforcement learning algorithms for the 3D tunnel game, certain aspects of the game learning step and parameter implementation are unique to this project and merit discussion. For instance, we

should look into the learning step in the game. The agent will not request a new move until the state has changed, despite the fact that it may seem more natural for a new decision to be made on every game tick. This has the effect of reducing the number of decisions that the agent must make in a given episode, but also results in intriguing policy behaviours that are elaborated upon in Chapter 6.

Furthermore, there are several parameters that are common to all learning agents, some of which were briefly discussed in the preceding chapter. In this section, we will examine in more detail how these parameters were integrated into this particular project. The initial optimistic value parameter is the simplest one to explain, as it is implemented in a straightforward manner. In particular, each time a state-action pair is added to the policy, its value is set to a predetermined number. This number (`initOptVal`), along with the other agent's sub-options parameters mentioned subsequently in this subsection, are specified through the command line (see 3.2).

The following two parameters worth noting are `eps` and `epsFinal`, which are responsible for the random moves executed by the ϵ -greedy policy. These parameters allow the user to specify the starting and ending values of ϵ . Then, at the end of each game, the new ϵ value is calculated by multiplying the current ϵ value with the *decrease* which is computed as follows¹:

$$decrease = \left(\frac{epsFinal}{eps} \right)^{\frac{1.0}{n}}$$

Here, `n` represents the number of games being played. The reason behind this epsilon decrease is to change the ratio between exploration and exploitation over time. At the beginning of the experiment, the `eps` value is higher, and thus random moves happen more often, causing the agent to try actions it would otherwise oversee. Later, when the policy is a bit stabilized, the `eps` value becomes smaller so it would allow the agent to play longer games and possibly win (if the `eps` value was high throughout the whole experiment, the agent would have a bigger chance of choosing an inadequate move and thus untimely ending the game).

Finally, it is pertinent to discuss the discounting value gamma (defined by `gam`). In this project, discounting is used in the following manner:

$$\gamma^{next_step.time - curr_step.time}$$

Normally, the gamma value is multiplied by itself on each step. However, as previously stated, learning steps in this implementation do not occur on every tick, but instead occur when the state changes. As a result, they may vary in size. To avoid uneven discounting, the time is calculated for each new decision made by the agent using the formula provided:

$$(game.num_of_ticks * 33) / 1000.0$$

5.3.2 Monte Carlo Agent

The Monte Carlo method is a type of reinforcement learning algorithm that updates its policy only after an episode is completed. This is done by iterating

¹This computation is performed once, when initializing the agent.

Algorithm 7 Updating policy for Monte Carlo

```
1:  $R \leftarrow \text{next\_step.score} - \text{curr\_step.score}$ 
2:  $G \leftarrow \text{POW}(\text{GAMMA}, \text{next\_step.time} - \text{curr\_step.time}) \cdot (R + G)$ 
3: if IS_FIRST_OCCURRENCE(...) then
4:    $\text{total\_return}[\text{curr\_step.state.action}] \leftarrow \text{total\_return}[\text{curr\_step.state.action}] + G$ 
5:    $\text{visits}[\text{curr\_step.state.action}] \leftarrow \text{visits}[\text{curr\_step.state.action}] + 1$ 
6: end if
```

through the entire episode, going from the last step towards the first, and increasing the number of visits and total return for each state-action pair, if this is their first visit inside this episode. The total return is calculated using the formula shown in Algorithm 7, while the number of visits is simply incremented by 1. To determine the optimal action, the agent compares the ratio of total return to number of visits for each possible action at a given state². This calculation is performed at each state transition during the episode. To clarify, instead of calculating a new move each time the `move()` function is called, the agents will always choose the same action based on the current state. Only once the state has changed, the new action is chosen based on the accumulated score and the new state. This implementation has resulted in a certain behaviour of the agents which will be more discussed in Chapter 6.3. As previously mentioned, the γ variable in the equation shown in 7 serves as a discount factor, meaning that the last move made, which resulted in termination of the game, will receive the highest penalty. As we move further down the list of moves, their significance decreases. It is important to note that if the value of γ is set to 1, all moves are given equal weight.

5.3.3 TD Agent

Unlike the Monte Carlo methods, which update their policies only after the completion of an episode, TD agents update their policies in real time, after each action is taken. To accomplish this, all TD agents have a shared function called `move()`, which calls the function displayed in 8.

Algorithm 8 Updating policy for TD Agent

```
1:  $\alpha \leftarrow 1.0 / \text{visits}[\text{last\_state.action}]$ 
2:  $\text{new\_state\_val} \leftarrow \text{new\_state.action}$ 
3: if terminal then
4:    $\text{new\_state\_val} \leftarrow 0$ 
5: end if
6:  $\text{new\_gamma} \leftarrow \text{POW}(\text{GAMMA}, \text{curr.time} - \text{prev.time})$ 
7:  $q[\text{last\_state.action}] \leftarrow q[\text{last\_state.action}] + \alpha(\text{new\_gamma}(R + \text{new\_state\_val}) - q[\text{last\_state.action}])$ 
```

²The pseudocode in Algorithm 1 a list of all returns for each state/action pair is kept, while in our implementation a return for a particular state/action pair is calculated with the mentioned ratio.

The update of the policy for a specific state value in TD learning involves multiple variables, most of which have been previously discussed. One new variable is α , which represents one divided by the total number of visits for a given state action pair. Furthermore, to make this calculation, a variable uniquely computed by each TD algorithm, known as `new_state_action`, is required. Various methods for computing this variable can be found in Algorithms 9, 10 and 11. If the current state is a terminal state, the `new_state_action` will not be used and instead, it will be replaced with a value of 0, as indicated in the update code.

Algorithm 9 Update function for SARSA

```

1: function GET_UPDATE(state, new_action, best_action)
2:   return  $q[\text{get\_state\_action}(\textit{state}, \textit{new\_action})]$ 
3: end function

```

Algorithm 10 Update function for Q-Learning

```

1: function GET_UPDATE(state, new_action, best_action)
2:   return  $q[\text{GET\_STATE\_ACTION}(\textit{state}, \textit{best\_action})]$ 
3: end function

```

Algorithm 11 Update function for ExpectedSARSA

```

1: function GET_UPDATE(state, new_action, best_action)
2:    $sum \leftarrow 0.0$ 
3:   for action in ACTIONS do
4:      $probability \leftarrow \frac{EPSILON}{\text{len}(\textit{ACTIONS})}$ 
5:     if action = best_action then
6:        $probability \leftarrow probability + 1 - EPSILON$ 
7:     end if
8:      $sum \leftarrow sum + probability \cdot Q(\textit{state}, \textit{action})$ 
9:   end for
10:  return  $sum$ 
11: end function

```

In SARSA, the `new_state_val` is calculated based on the value of the next action the agent will take, denoted as `new_action`. On the other hand, Q-learning uses the value of the best action possible in the next state, denoted as `best_action`. These two variables are equal if a greedy policy is implemented. However if we consider a ϵ -greedy policy, then they might differ based on whether a random action has been chosen. Expected SARSA combines these two approaches by taking the expected value of all possible actions in the next state.

Similar to the agents in the TD class, the update for the Double Q-Learning agent occurs each time the agent changes its state. The update process is slightly different. In this method, two separate action-value functions, denoted as Q1 and Q2, are used to estimate the maximum action value for a given state. At each update step, one of the Q-values is selected randomly and updated using the other Q-value as a reference. This process helps to reduce the overestimation of action values and leads to more stable learning.

Algorithm 12 Update function for Double Q-Learning

```
alpha  $\leftarrow$  1.0/visits[last_state_action]  
new_gamma  $\leftarrow$  POW(GAMMA,next_step.time - curr_step.time)  
if rand.randf_range(0, 1) < 0.5 then  
    new_state_val  $\leftarrow$  q2[get_state_action(state, best_action)]  
    if terminal then:  
        new_state_val  $\leftarrow$  0  
    end if  
    q[last_state_action]  $\leftarrow$  alpha(new_gamma(R + new_state_val) -  
    q[last_state_action])  
else  
    new_state_val  $\leftarrow$  q[get_state_action(state, best_action)]  
    if terminal then:  
        new_state_val  $\leftarrow$  0  
    end if  
    q2[last_state_action]  $\leftarrow$  alpha(new_gamma(R + new_state_val) -  
    q2[last_state_action])  
end if
```

6. Experiments

In this chapter, we will evaluate the performance of the various agents when confronted with different combinations of obstacles, as well as presenting interesting observations made during the experiments. One of the key questions we aim to answer is whether any of the agents are capable of learning to play the entire game. To begin, it is important to explain how the experiments were conducted and the significance of the plots in the following subsections of this chapter.

6.1 Conducting the experiments

The process of conducting experiments encountered several challenges and fluctuations. Initially, the primary concern was to determine the appropriate number of rotations (**rots**) and distances (**dist**s) for each obstacle type (explained in Table 3.2). After experimenting with up to 30 rotations in some cases, and not getting satisfying results with seemingly any combination of other parameters, it was determined that the game’s difficulty in later levels was the root of the problem, as confirmed by human players. Consequently, the game had to be adjusted and thus, there are slight variations in parameters, such as the starting speed and distances between obstacles, in the version of the “Space-run” game used in this thesis, as compared to the original. As per a human player’s assessment, it is now possible to play all environment combinations until level 15 or even beyond. It is noteworthy that level 10 was the initial choice for the winning level, which was later shifted to level 15 to prevent the agent from settling for a mediocre policy and to find the optimal policy. However, this shift did not yield significant results, and the same behaviour could likely be achieved by increasing the number of games, allowing the ϵ -greedy policy to perform random moves more frequently for an extended period. Despite this, level 15 was used in all subsequent experiments.

LadybugFlying	⇒ 6	Rotavirus	⇒ 6	MovingI	⇒ 6	O	⇒ 13
LadybugWalking	⇒ 6	Tokens	⇒ 6	HexO	⇒ 7	Balls	⇒ 15
Worm	⇒ 6	HalfHex	⇒ 6	Triangles	⇒ 7	Hex	⇒ 22
Bacteriophage	⇒ 6	I	⇒ 6	X	⇒ 11	Walls	⇒ 22

Figure 6.1: **rots** values used for each obstacle type

After addressing the the issue of game not being playable, the most straightforward method for identifying the number of rotations required was to play the game manually in debug mode. The results obtained from these experiments are displayed in Figure 6.1 and were used in all experiments conducted. However, for obstacles such as **Walls** and **Balls**, this method was not viable. The reason being that running `env=Balls` or `env=Walls` with visuals caused the game to lag significantly due to the number of animations playing simultaneously. As a result, **Walls** received the same number of rotations as **Hex** trap, while **Balls** received 15 rotations, the number at which the agent managed to learn. For bug, virus,

and token obstacles, the default and minimum value of 6 rotations was assigned, with which they all trained successfully. Concerning the `dists` parameter, it was concluded during the experiments that all agents could learn any obstacle with `dists=1`, and increasing this number needlessly would only increase the number of states required for training.

The remaining values to be determined for the experiments were the suboptions for each agent (see 3.2.1). The reader must realize that there were a lot of experiments conducted that had to be disregarded later on for this or that reason, but they provided useful insights that could be utilized. For instance, in many of those experiments, the agent performed best with `eps` values ranging from 0.2 to 0.4, with 0.2 being the most common, in combination with an `epsFinal` value of 0.0001. The rationale for using a low `epsFinal` value is that towards the end, the agent almost exclusively exploits the current policy, and a more gradual decrease in epsilon values is suitable for larger values of `n`. Furthermore, `initOptVal` of 20.0 and 100.0 was promising in most experiments. The `gam` value will be discussed in a later subsection.

Throughout the months dedicated to conducting experiments, they gradually converged towards checking combinations of the values mentioned above. In some instances, other values were experimented with and will be discussed when exploring various environments in later subsections. However, due to the length of the experiments and the number of parameters requiring adjustment, they were kept systematic towards the end, and most of the results presented in this chapter were obtained from experiments run with different `env` and value `n` adjusted according to the environment. Each combination was tested on ten seeds for each agent, with the combinations consisting of `eps` parameter taking a value of either 0.2 or 0.4, and the `initOptVal` being either 20.0 or 100.0. Discounting was kept at the value of 1 and `epsFinal` was 0.0001.

6.2 Plotting

In this chapter, a lot of plots similar to the one shown in Figure 6.2 will be presented. A subsection has been dedicated to explaining the different components of these plots and how to read them. While some plots may differ slightly from the one in Figure 6.2, the meaning behind them can still be easily deduced.

The top part of the figure displays all the necessary information that was used in the experiment. Most of the values have been previously described, except for “previous game”. This value is meant for experiments that used the `database=read` option and were performed on an agent that had trained for some number of games. This number specifies the how many of games the agent had previously trained on.

Moving further down, there is a plot with three lines and a legend in the top right corner. The data line, as indicated on the figure, represents the score that the agent achieved on a particular episode. However, some plots may show the average value of the agent’s score for that episode over several different seeds used for the random action. Additionally, there may be multiple data lines on the plot, each representing a different agent. In these cases, each line is averaged between many seeds. Each agent is labelled with their respective color inside the legend.

The mean line represents the average score value for the entire experiment,

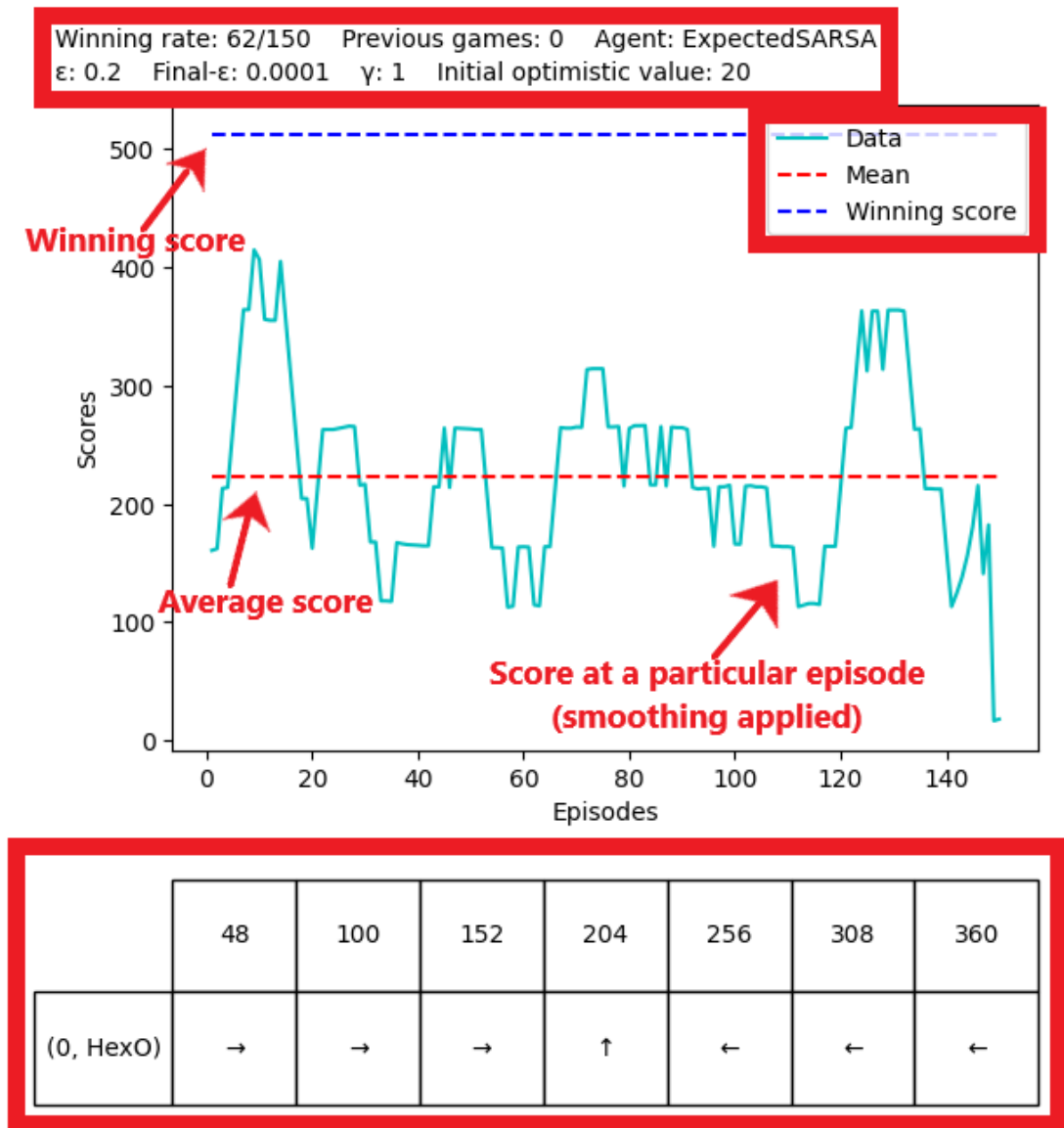


Figure 6.2: Plot example

while the winning score represents the winning threshold. It should be noted that if the agent did not win any games, or the plot is displaying multiple agents, the winning score line will be omitted. Furthermore, the data line is averaged to appear smoother. This is why even though the winning rate is 62/150, the data line doesn't touch the winning score threshold at any point.

As mentioned in the Section 3.2, it is possible to perform continuous evaluation on experiments, meaning one learning game is played with random actions, followed immediately by an evaluation game using only the current policy that the agent is performing. In all experiments conducted, the data line only represents the evaluation games.

At the very bottom of the figure, there is a table which only appears in plots that contain a single data line that is not averaged over different seeds and has only one seed value. In the table, all rotation values for this experiment define the columns, while each row has a tuple of distance value and type of the obstacle. Since all experiments are performed on `dist=1`, the number of rows will match

only the number of different obstacles used in the experiment.

It should be emphasized that parameters such as the number of seeds in the averaged data line or the size of the smoothing window will be clearly specified for each plot mentioned in the rest of the chapter. This prevents any form of ambiguity or confusion regarding the details of the experiment.

6.3 Interesting behaviours

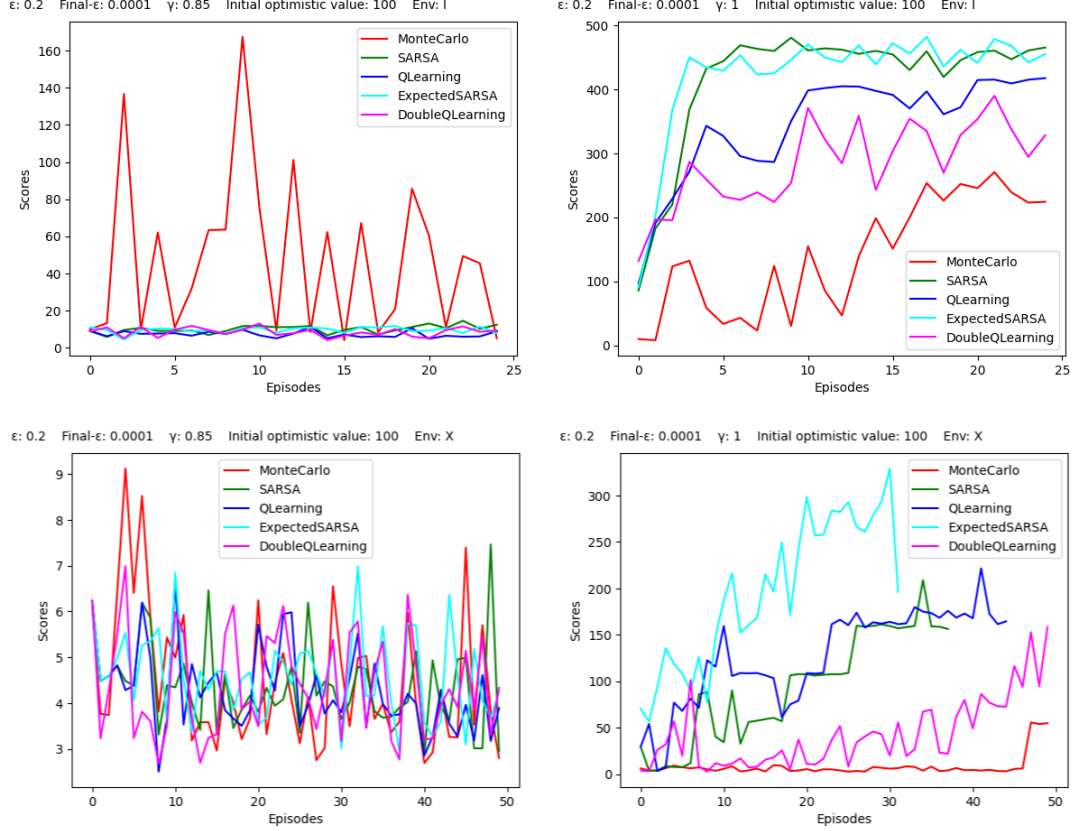


Figure 6.3: Discounting example

In this chapter, we aim to discuss certain unexpected findings that surfaced during our experimentation. One of the immediate observations can be seen in Figure 6.3¹. We conducted experiments on two distinct environments, **env**=[I] and **env**=[X], and for each environment, we carried out experiments with **gam**=0.85 and **gam**=1.0 for all agents. As evident from the plots, the lower gamma value exhibited considerably poorer performance than when no discounting (**gam**=1.0) was applied. This trend is not limited to these specific environments and testing conditions but rather observed consistently across all our experimentation. This result is counter-intuitive since it seems logical that penalizing the last action more than previous ones would result in a better policy.

Upon further investigation, we discovered that in some cases, a lost game for the agent does not result from the last action directly but rather from a

¹No smoothing was applied to any of the plots in this subsection

```

epsilon = 0.0764
  60 120 180 240 300 360
(0, I) ^ < > > >
evaluation game:
died on level 13, rot = 360
Game 38 score: 451.7
last actions: [0,60,I]_[0,0] [0,360,I]_[1,0] [0,60,I]_[0,0] [0,180,I]_[1,0] [0,240,I]_[1,0] [0,300,I]_[1,0] [0,360,I]_[1,0]

```

Figure 6.4: Discounting explanation

chain reaction initiated by a previous “bad decision”. As seen in Figure 6.4, the agent’s last action of going right at rotation 360 to reach a safe one, 60, is not a poor decision in itself but rather the best possible action in that state². However, analysing the last four actions taken by the agent, it becomes clear that it attempted to reach rotation 60 by going left from the rotation 180. With a high score of 451.7 (Figure 6.3), the agent was undeniably fast, and while this policy may have been effective in the early game before the agent attained its current speed³, there is simply not enough time for the agent to rotate at this point in the game. In this case, one could argue that the fourth action from the last was responsible for the agent’s loss. For cases like this, we believe that the agent performs better when all actions are penalized equally.

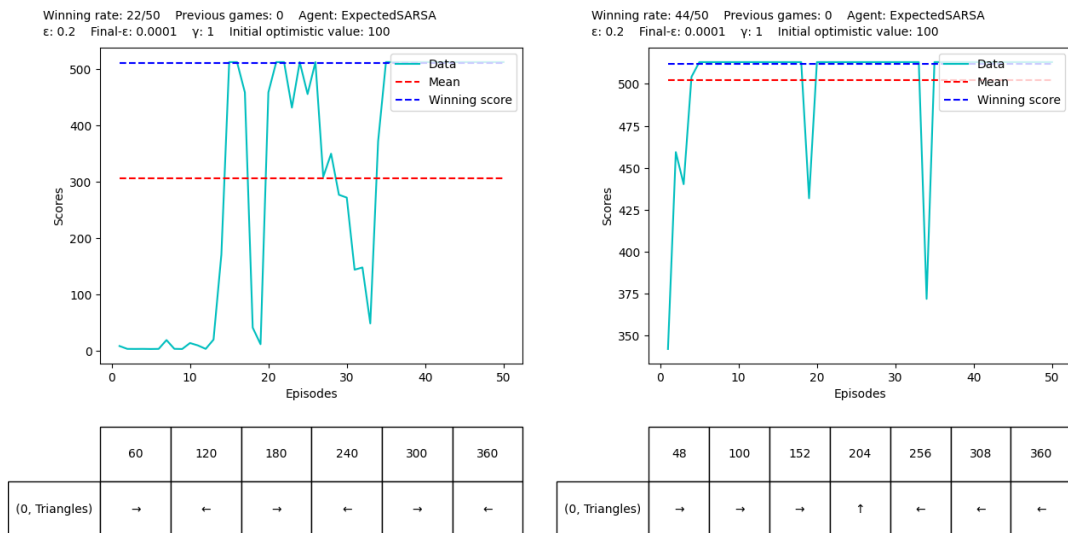


Figure 6.5: Triangles example

Another intriguing behaviour emerged as a result of an unforeseen coincidence. As previously noted, we chose to have the agent not take a new action every time its `move()` function is called, but rather return the same action until the state changes. This resulted in a substantial reduction in the number of different actions taken by the agent, given that the move function is invoked every game tick, while state changes occur less frequently. This approach led to a behaviour that could not have been predicted at the outset. Figure 6.12 provides an illustration of this behaviour (with some sentiment `env=[Triangles]` was chosen as this was the first environment on which the behaviour was noticed).

²As confirmed by a human player, rotation 60 is safe for trap type I

³It should be recalled that after every 3 levels the agent’s speed increases

Ordinarily, `env=[Triangles]` does not offer any safe rotations unless `rots=7` or more. As shown in the plot, the agent will certainly learn with this rotation value. However, if the agent is given only 6 rotation values to choose from, it will develop a policy that rapidly oscillates between two rotations, keeping the player character, Hans, on the edge of those rotations, allowing him to safely pass through the trap. This behaviour is not confined to situations where the agent is “forced” to make such a decision. During training, in many instances, the agent will learn to stay on the edge rather than advance into a completely safe rotation. There does not appear to be a preference for one or the other; rather, the policy the agent discovers first is determined by other experimental factors. It can be concluded that this behaviour arose solely because the agent did not alter its action until the state changed. In my opinion, this discovery is one of the most exciting outcomes of this project.

6.4 Individual traps environment

This section discusses the performance of agents in an environment containing only one trap at a time. Considering the number of traps and agents, for each trap, we have picked a plot that showed the best performances of the agents overall. Some specific cases will, of course, be highlighted. It should be noted that all plots within these examples have had smoothing applied with `window=10`. Thus certain spikes will not be visible. Additionally, unless otherwise suggested, you can assume the values that were produced are an average of 5 different seeds. The aim was to show a realistic picture on how the agent would perform, and not only highlight the occurrences in which the outcome was satisfactory but rather account for the failures in reproducing the perfect policy as well.

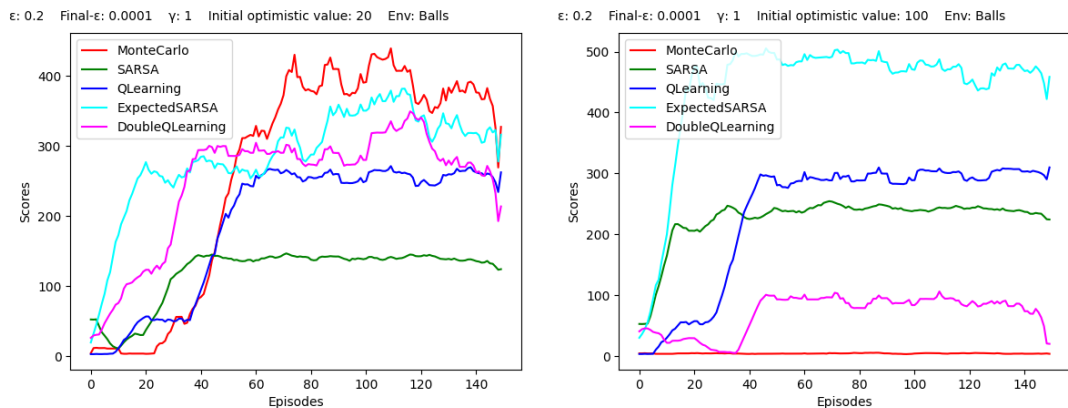


Figure 6.6: Balls trap examples

Another very important factor were the choices of the hyper parameters. As mentioned before, most commonly used one will be `eps=0.2` or `ep=0.4` and `initOptVal=20.0` or `initOptVal=100.0`. Perhaps best visible in Figures 6.6 and 6.7, is how much hyper parameters can influence the outcome. Both the Balls and Hex0 traps varied significantly based on the initial optimistic value used, with the other agents’ performances also varying to a lesser extent. These

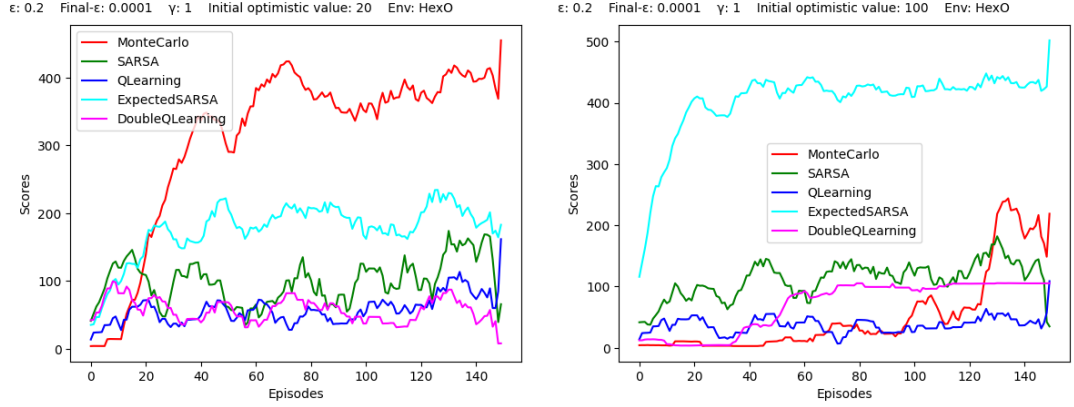


Figure 6.7: HexO trap examples

findings underscore the importance of carefully selecting hyper parameters when training reinforcement learning.

Nevertheless, in some cases we have not managed to find the right combination of the hyper parameters. The Hex trap is a prime example of such a challenge, requiring no less than 22 rotations to be navigated successfully. As shown in Figure 6.8, the agents' overall performance on this trap has been moderate at best.

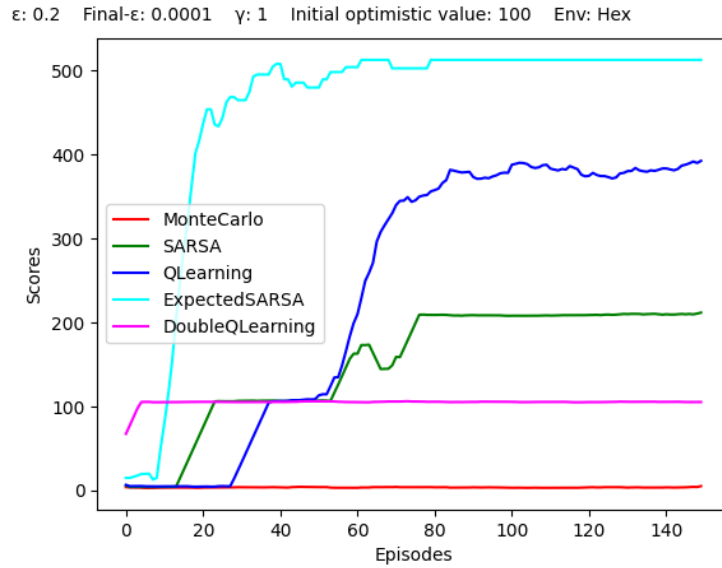


Figure 6.8: Hex trap example

However, Figure 6.9 presents the results of an experiment conducted on the **MonteCarlo**(MC) agent, which was produced by a single seed value. The choice of seed or hyper parameter values is not critical, as all the values for the Hex trap for this particular agent appear the same. By examining the policy that the MC agent learned in this experiment and comparing it to the left-hand plot, which depicts the **DoubleQLearning**(DQL) agent learning an arguably perfect policy right from the start, we can see that the MC agent completely overlooks the safe rotations and instead chooses to move forward in areas that lead to either certain death,

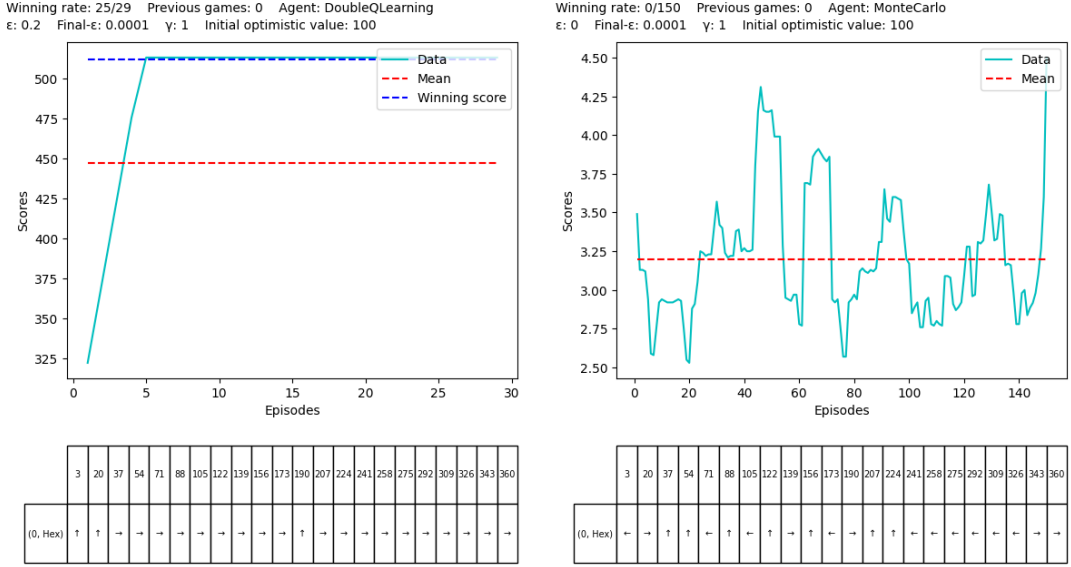


Figure 6.9: Hex trap with DoubleQLearning and MonteCarlo agents

or half-safe rotations which let the agent pass based on pure chance. It is worth noting that the performance of the DQL agent varies depending on the seed value (on the plot we can see that the average performance is even lower than with other agents), and this example was included only to illustrate what an optimal policy for this type of trap might look like.

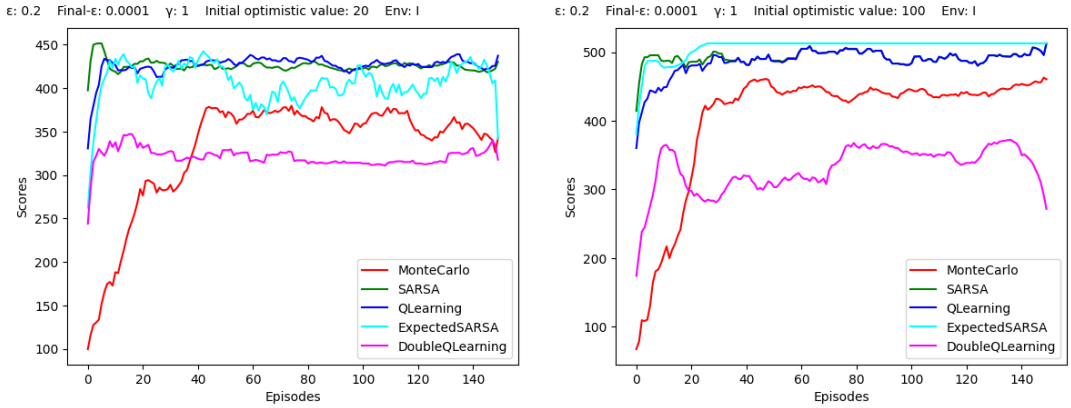


Figure 6.10: I trap example

It is worth noting that there exist cases where the aforementioned variations in hyper parameter selection lead to highly desirable outcomes. This is exemplified by the results presented in Figures 6.10 and 6.11, which demonstrate the efficacy of the ExpectedSARSA(ES) agent. Notably, on the right hand side of the both plots, the ES agent was able to achieve optimal performance early on in the game with all seed values. ES agent, in a very large number of experiments, has outperformed its counterparts, sometimes by a significant amount. This is particularly evident when considering its performance in simpler environments such as traps I and MovingI, where it is apparent that the agent is capable of learning extremely well. While other agents have performed well on these specific

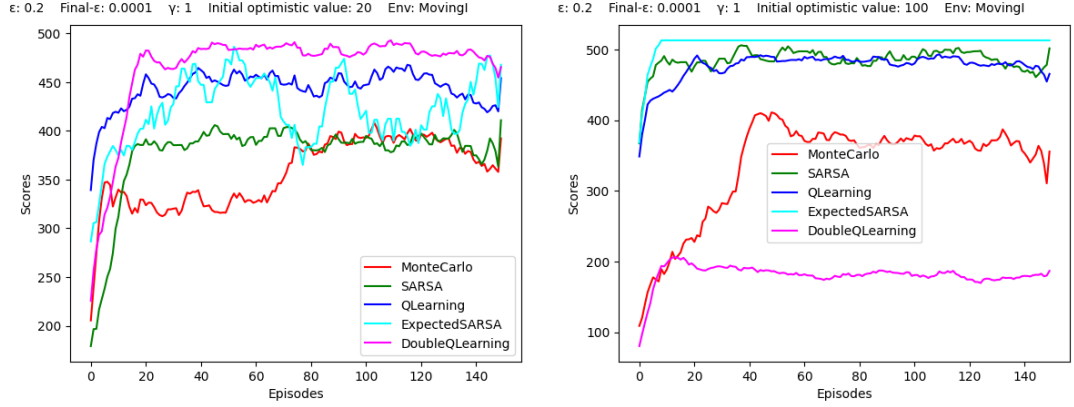


Figure 6.11: MovingI trap example

traps as well, their success may not be immediately apparent from the averaged results depicted in the plots.

As can be observed from some of the previous graphs, agents often exhibit significant changes in behaviour with different hyper parameters. However, as depicted in Figure 6.12 for the **Triangles** trap, increasing the `initOptVal` resulted in the agents maintaining the same ratio between each other, despite three of them increasing their average score significantly.

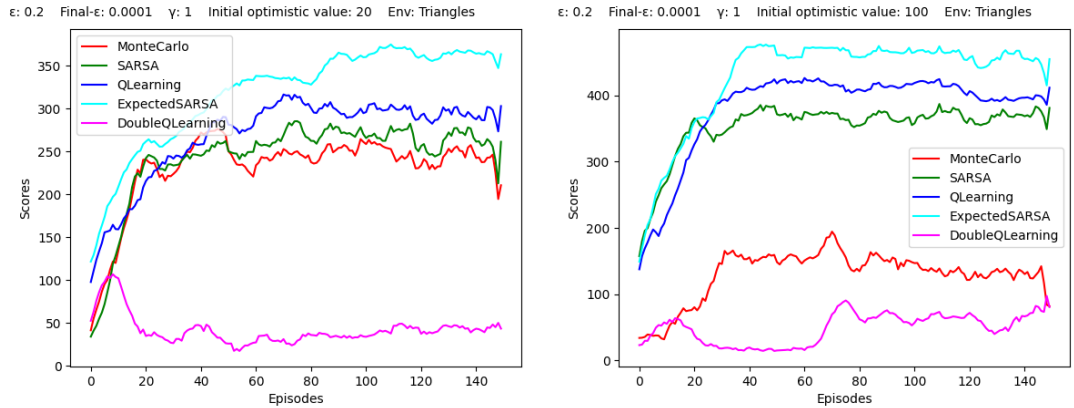


Figure 6.12: Triangles trap example

Figure 6.13 showcases the performance of the agents on the remaining traps that have not been mentioned before in this section. It is worth noting that the agents' overall behaviour on these traps is satisfactory, with some agents performing exceedingly well. This, once again, reinforces the idea that the choice of hyper parameters can have a significant impact on the agent's performance, and that a careful selection of these parameters is crucial for achieving optimal results.

Although performing these experiments with a larger number of seeds would ideally yield even more accurate results, we hope that the picture we presented provides a reasonable representation of the agents' performance.

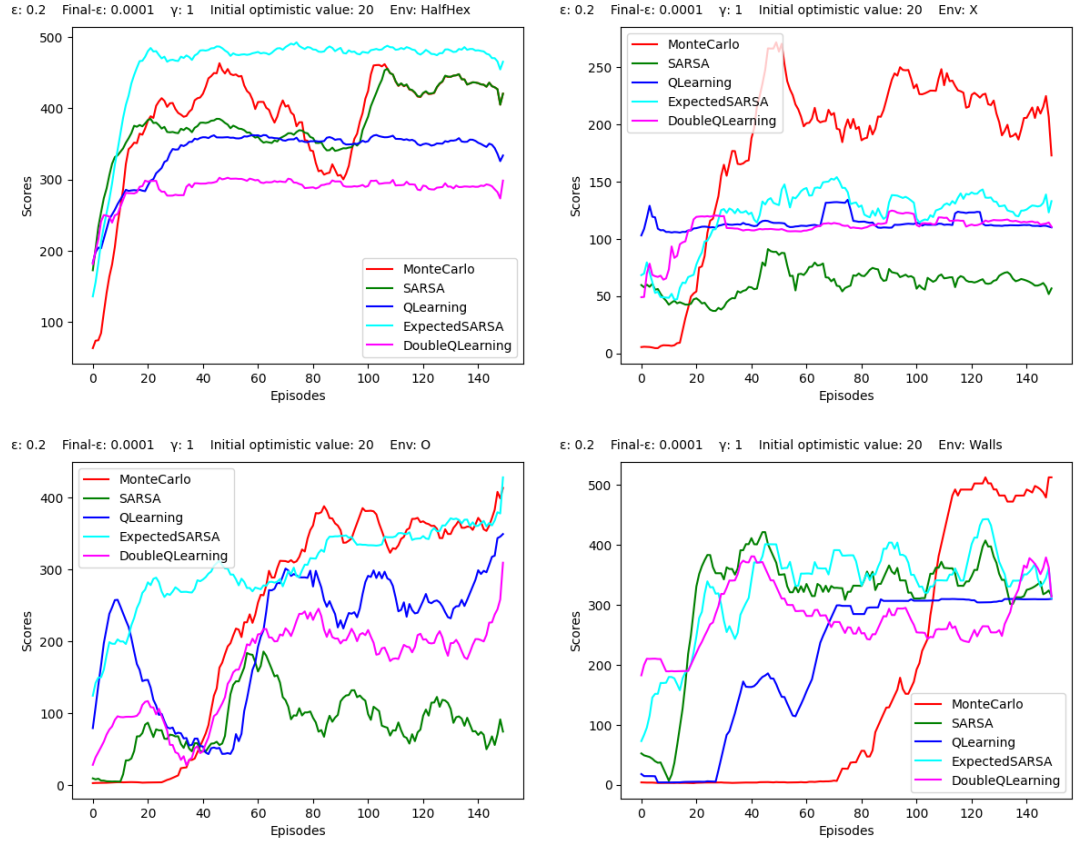


Figure 6.13: HalfHex, X, O and Walls trap examples

6.5 Traps environment

6.6 Tokens environment

6.7 Bugs environment

6.8 Viruses environment

6.9 Full game environment

Conclusion

In conclusion, this thesis has examined the performance of different reinforcement learning algorithms on various environments in the Godot game engine. The results showed that the Expected SARSA algorithms performed consistently well in multiple environments, while the Double Q-Learning algorithm struggled with some but improved with more training. All algorithms performed quite well with the individual traps, while with all trap types combined, Monte Carlo and once again Expected SARSA gave the best results. In the token environment, all algorithms performed well. A very interesting observation is that in all cases that include bugs and viruses, the agents chose to rather avoid the obstacle than shoot it down. Unfortunately, none of the algorithms were able to achieve optimal performance on the full game. Further research is needed to determine the optimal parameters and strategies for these algorithms to accomplish that task.

Bibliography

Tunnel Rush, 2023. URL <https://tunnelrush2.com/>.

Una Adilović. Space Run AI, 2023. URL <https://github.com/AdilovicUna/Space-run-AI>.

Khronos. glTF - Runtime 3D Asset Delivery, 2023. URL <https://www.khronos.org/glTF/>.

Juan Linietsky. Godot Engine - Documentation, 2021. URL <https://docs.godotengine.org/en/stable/index.html>.

Ton Roosendaal. Blender, 2023. URL <https://www.blender.org/>.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, Massachusetts / London, England, second edition, 2018.

List of Figures

1.1	Movement	4
1.2	Hans	5
1.3	Trap examples	5
1.4	Bugs	5
1.5	Bacteriophage and Rotavirus	6
1.6	Battery and Energy Token	7
2.1	Structure of Game.tscn	8
2.2	State	11
3.1	Rotations when rots = 6	14
5.1	Agents hierarchy inside the project	26
6.1	rots values used for each obstacle type	32
6.2	Plot example	34
6.3	Discounting example	35
6.4	Discounting explanation	36
6.5	Triangles example	36
6.6	Balls trap examples	37
6.7	HexO trap examples	38
6.8	Hex trap example	38
6.9	Hex trap with DoubleQLearning and MonteCarlo agents	39
6.10	I trap example	39
6.11	MovingI trap example	40
6.12	Triangles trap example	40
6.13	HalfHex, X, O and Walls trap examples	41

List of Tables

List of Abbreviations

A. Attachments

A.1 First Attachment