

UNIVERSITY OF SÃO PAULO  
INSTITUTE OF MATHEMATICS AND STATISTICS  
BACHELOR OF APPLIED MATHEMATICS

**Simulation of a self-parking car using  
deep reinforcement learning**

André Shumhei Kato

FINAL ESSAY  
MAP 2010 — CAPSTONE PROJECT

Supervisor: Prof. Dr. Alexandre Megiorin Roma

São Paulo  
2022

*The content of this work is published under the CC BY 4.0 license  
(Creative Commons Attribution 4.0 International License)*

*Dedico esse trabalho ao meu pai e maior mentor,  
Marco, pelo apoio incondicional durante toda a  
minha trajetória e por nunca ter medido esforços  
para me ajudar, me permitindo seguir meus sonhos.*



# Agradecimentos

*Fall seven times, get up eight.*

— Japanese Proverb

Ao meu orientador Alexandre Roma por toda a sua atenção, cuidado e orientação.

Aos meus amigos Vinícius Moreno e Lucas Dário pelas suas contribuições ao longo da realização desse trabalho

Ao meu amigo Kauer Kitinger d’Avila por ajudar na revisão da monografia.

Aos meus amigos Victor, Luma, João, Yukio, Andrew, Carlos Wellington e muitos outros que gostaria de adicionar à já extensa lista, pelos momentos de diversão e desconpressão.

Aos professores David Silver da University College London e Emma Brunskill de Stanford por disponibilizar gratuitamente aulas e materiais de estudo.

Ao **Samuel Arzt** da University of Applied Sciences Salzburg, autor do **vídeo** que motivou e inspirou a realização desse trabalho.



# Resumo

André Shumhei Kato. **Simulation of a self-parking car using deep reinforcement learning**. Monografia (Bacharelado). Instituto de Matemática e Estatística, Universidade de São Paulo, São Paulo, 2022.

Neste trabalho, apresentamos uma visão geral sobre o aprendizado por reforço, com foco em *proximal policy optimization* (PPO). Nos últimos anos, muitos métodos em aprendizado por reforço surgiram, incluindo os baseados em *policy gradient* (PG). Os métodos PG têm algumas vantagens relevantes sobre outros métodos, uma vez que têm boas propriedades de convergência e lidam bem com espaços de ação contínuos. Estamos particularmente interessados no estudo do PPO, um método de PG derivado de *trust region policy optimization* (TRPO), que foi proposto como uma melhoria em relação ao PG original. Ambos os métodos previnem atualizações de peso destrutivas, que colapsam o processo de treinamento. PPO tem uma ligeira vantagem sobre a TRPO por ser um método menos complexo de derivar e menos computacionalmente complexo em comparação ao TRPO. Combinando PPO com métodos de aprendizagem profunda, estudamos um problema de estacionamento, onde temos um carro, cujo objetivo é estacionar em uma vaga designada em diferentes circunstâncias. Também discutimos o *design* das recompensas nos experimentos, mostrando alguns dos problemas encontradas durante o desenvolvimento e a forma como foram resolvidos.

**Palavras-chave:** Machine learning. Unity. Inteligência artificial. Redes neurais. Aprendizado por reforço. Proximal policy optimization. Estacionamento automático.





# Abstract

André Shumhei Kato. **Simulation of a self-parking car using deep reinforcement learning**. Capstone Project Report (Bachelor). Institute of Mathematics and Statistics, University of São Paulo, São Paulo, 2022.

In this work, we present a general overview on deep reinforcement learning (RL), focusing particularly on proximal policy optimization. Over the recent years, many RL methods have been derived, including ones based on policy gradient (PG). Policy gradient methods have important advantages over other methods, since it has good convergence properties and can handle continuous action spaces well. We are particularly interested in proximal policy optimization (PPO), a policy gradient method derived from trust region policy optimization (TRPO), which was proposed as an improvement over vanilla policy gradient. Both of these methods avoid destructive weight updates, a issue commonly found on vanilla PG, with PPO having a slight edge on TRPO for being less complex to derive and less computationally expensive compared to TRPO. Combining PPO with deep learning methods, we tackle a car parking problem, where a car has to park in a designated spot under different circumstances. We also discuss the reward design for the experiments, showing some of the pitfalls encountered during development and how they were addressed.

**Keywords:** Machine learning. Unity. Artificial intelligence. Neural networks. Reinforcement learning. Proximal policy optimization. Autonomous parking.



# Lista de Abreviaturas

ML	Machine Learning
RL	Reinforcement Learning
MDP	Markov Decision Process
PPO	Proximal Policy Optimization
TRPO	Trust Region Policy Optimization
PG	Policy Gradient
ANN	Artificial Neural Network
NN	Neural Network
FF	Feed-forward
MLP	Multilayer Perceptron
ReLU	Rectified Linear Unit
ELU	Exponential Linear Unit
GLP	Grad-Log-Prob
EGLP	Expected Grad-Log-Prob
IME	Instituto de Matemática e Estatística
USP	Universidade de São Paulo



# Contents

<b>1</b>	<b>Preliminaries</b>	<b>1</b>
1.1	Introduction . . . . .	1
1.2	Motivation . . . . .	1
<b>2</b>	<b>Introduction</b>	<b>3</b>
2.1	Reinforcement Learning . . . . .	3
2.2	Elements of a Reinforcement Learning Problem . . . . .	4
2.3	Agent-Environment Interface . . . . .	5
2.4	Goals and Rewards . . . . .	5
<b>3</b>	<b>Markov Decision Processes</b>	<b>7</b>
3.1	Markov Processes . . . . .	7
3.2	Policies . . . . .	8
3.3	Optimality . . . . .	9
3.4	Designing Rewards . . . . .	10
<b>4</b>	<b>Neural Networks</b>	<b>13</b>
4.1	Machine Learning . . . . .	13
4.2	Neural Networks . . . . .	14
4.3	Perceptron . . . . .	15
4.4	Multilayer Perceptron . . . . .	17
<b>5</b>	<b>Proximal Policy Optimization</b>	<b>21</b>
5.1	Policy Gradient . . . . .	21
5.2	Trust Region Policy Optimization . . . . .	24
5.3	Proximal Policy Optimization . . . . .	25
<b>6</b>	<b>Experiments</b>	<b>27</b>
6.1	Hyperparameters and Evaluation Metrics . . . . .	27

6.2	Experiment 1: Fixed Positions . . . . .	28
6.3	Experiment 2: Randomized Car Position and Fixed Parking Spot . . . . .	33
7	<b>Conclusions</b>	<b>35</b>
	 <b>References</b>	 <b>37</b>

# Chapter 1

## Preliminaries

### 1.1 Introduction

Reinforcement Learning is considered a subfield of Machine Learning, where learning occurs through an agent interacting with an environment. At each time step, the agent performs an action and the environment responds by producing a reward signal and transitioning to the next state. The goal of the agent is to maximize the total expected reward (Richard S. SUTTON and BARTO, 2015).

There are a lot of challenges that naturally arise from reinforcement learning problems that differ from the ones faced in classical machine learning. For example, balancing immediate rewards and future rewards: up to which point is it worth to sacrifice early rewards in exchange for bigger rewards in the future? An agent might be inclined to take actions that have been taken before because it has the knowledge of how much reward those actions will yield. Thus, limiting the agent to that specific set of actions and ultimately impairing it from further exploring the environment and possibly discovering new states and actions that could yield even more rewards. This is called the exploration-exploitation dilemma (Richard S. SUTTON and BARTO, 2015).

The focus of this work will be an application of the Proximal Policy Optimization (PPO) algorithm to train an agent able to park a car in a designated spot.

### 1.2 Motivation

With car crashes being more and more common, car manufacturers have started working on technologies that help mitigate crashes, ranging from simple proximity sensors that alarms the driver if a collision is imminent or stops the car by itself to fully fledged auto-driving systems. In the latter, automated parking is a key part in autonomous vehicle systems that allows cars to navigate through a parking lot completely unassisted.

This work aims to recreate the self-parking system inside a 3D virtual environment using deep reinforcement learning and studying how the algorithm performs in different parking situations.





# Chapter 2

## Introduction

### 2.1 Reinforcement Learning

Reinforcement learning is learning what to do - mapping situations to actions in order to maximize the reward received. The learner (or "agent") has no knowledge on what actions should be taken, as in many forms of machine learning. Instead, it must discover which ones yield the most reward by trying them and observing what happens. When taking an action, that action may affect not only the immediate reward but also the next reward and all the subsequent others. These two characteristics - trial-and-error search and immediate vs. delayed rewards are the most distinguishing aspects of a reinforcement learning problem.

The following examples illustrate how reinforcement learning concepts are applied in real life:

1. A chess player making a move. The move is informed both by planning - anticipating possible replies and counterreplies - and by intuitive judgements of what positions and moves are desirable.
2. A gazelle calf struggling to stand on its feet after being born and a few hours later being able to run.

Both examples involve an interaction between a decision-making agent and the environment, in which the agent seeks to achieve a goal, despite uncertainty about its environment. The agent's action may affect the future state of the environment, for example, the next chess move will affect the possible options and opportunities in the future. Taking the correct choice requires taking into account indirect and delayed consequences of actions, and thus requires planning.

Both examples have goals in which the agent can judge the progress towards it based on what it can sense directly. For example, the chess player could judge his progress by comparing his remaining pieces with the opponent's. The player also knows whether he wins.

In order to fully formulate a reinforcement learning problem, we need optimal control

of a Markov Decision Process, which will be discussed later, but the basic elements are shown in the next section.

## 2.2 Elements of a Reinforcement Learning Problem

Apart from the agent and environment, there are other subelements of a reinforcement learning system: a *policy*, a *reward signal*, a *value function* and, optionally, a *model of the environment*.

A *policy* is a mapping from perceived states to actions to be taken when in those states, that is, a policy is what defines the agent's behavior. Policies can be deterministic, being a simple function or a lookup table, or stochastic, with probabilities associated with each action.

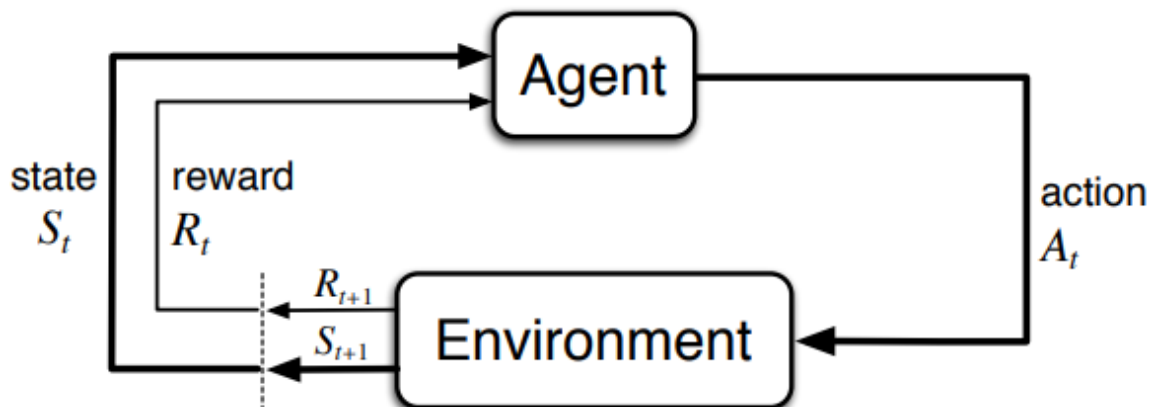
A *reward signal* is what defines the goal in a reinforcement learning problem. At each time step, the environment sends a reward signal, which is nothing but a number. The agent's goal is to maximize reward received over the long run. In biological systems, rewards are analogous to feeling pleasure or pain. The reward received depends on the agent's current action and on the state of the environment. The agent cannot alter this process in any way, but it can influence it through its actions. The reward signal is the primary basis for altering a policy. If an action selected by the policy yields a low reward, then the policy may be changed to select another action in that situation in the future.

A *value function* tells us how much reward the agent should expect to receive over the future by starting in a specific state. While reward signals indicates whether a state is immediately desirable or not, the value function estimates the long-term desirability of that state by taking into account the states that are likely to follow the current one. For example, a low reward state might be followed by a high reward state or vice versa. It's values we are most concerned with when making and evaluating decisions - we seek for actions that yield maximum value, not reward. Note that this is equivalent to maximizing reward over the long run. The major problem with value is that it can be hard to estimate - rewards are given directly by the environment, but values must be re-estimated at each iteration from sequences of observations an agent makes over its lifetime. The most important component of almost all reinforcement learning algorithms is value function estimation.

A *model of the environment* is something that mimics the behavior of the environment itself. More generally, that allows for inferences about how the environment will respond to certain actions. For example, for a given pair of state and action, the model might predict how the environment will respond to that action in that particular state. Models are particularly useful for *planning*, which is deciding on a course of action by considering possible future situations. Methods that use models and planning are called *model-based* methods, as opposed to *model-free* methods that rely on trial-and-error to learn about the environment.

## 2.3 Agent-Environment Interface

The learner and decision-maker agent and the thing it interacts with, comprising of everything outside the *agent*, is defined as the *environment*. The agent and the environment interact continually, with the agent choosing actions and the environment responding to those actions, presenting new situations to the agent and rewarding it.



**Figure 2.1:** Representation of the agent-environment interface from *Richard S. SUTTON and BARTO, 2015*.

In other words, at each time step  $t = 0, 1, 2, \dots$ , the agent receives some representation of the environment's state  $S_t \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of all possible states, and selects an action  $A_t \in \mathcal{A}(S_t)$ , where  $\mathcal{A}$  is the set of all possible actions and  $\mathcal{A}(S_t)$  the set of all possible actions in state  $S_t$ . At the next time step  $t + 1$ , the agent receives a reward  $R_{t+1} \in \mathbb{R}$  and a new state  $S_{t+1}$ .

The action the agent takes is sampled from the agent's policy, denoted  $\pi$ , with  $\pi(a | s)$  being the probability of taking action  $a$  when in state  $s$ .  $\pi$  is just an ordinary function defining a probability distribution over  $a \in \mathcal{A}(s)$ . Reinforcement learning methods specify how the policy is changed as the agent gathers more experience.

## 2.4 Goals and Rewards

The goal of an agent is formalized in terms of a reward signal  $R_t \in \mathbb{R}$  that is passed to the agent by the environment. Consider the sequence of rewards received after time step  $t$ :  $R_{t+1}, R_{t+2}, \dots$ . We define a *return*  $G_t$ , which is a function of that sequence. The simplest case is the sum of all of them:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

The agent's goal is to maximize the total amount of rewards it receives. In other words, the goal is to maximize the *expected return*  $\mathbb{E}[G_t]$ . In the above definition, we have the notion of a final time step, which comes very naturally when the agent-environment interaction can be broken down into subsequences, which we call *episodes* or *trials*. For example, plays

of a game, trips through a maze or any sort of repeated interactions can be considered episodes. Each episode ends in a special state called *terminal state*, followed by a reset to a standard state or to a sample of a distribution of starting states.

However, not all agent-environment interactions can be broken naturally into episodes, instead, it could go on without limit. That is, with slight abuse of notation,  $T = \infty$  and the return  $G$  could easily be infinite as well. Thus, the concept of *discounting* arises, where we introduce a *discounting factor*  $\gamma$ ,  $0 \leq \gamma \leq 1$ , to weigh immediate and future rewards:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Then, the sum is finite as long as  $\gamma < 1$  and the reward sequence is bounded. Also note that if  $\gamma = 0$ , the agent is only concerned with immediate rewards and the action  $A_t$  will be chosen in such a way to maximize only  $R_{t+1}$ . If  $\gamma$  is close to 1, then the future rewards will be taken into account more strongly.

According to [Richard S. SUTTON and BARTO, 2015](#), although formulating goals in terms of reward signals appears to be limiting, in practice, it has proved to be flexible and widely applicable. For example, to make a robot learn how to escape from a maze, the reward is often -1 for every time step that passes prior to escape; this encourages the agent to escape as quickly as possible.

While designing how rewards should be given to the agent, it's crucial that it's done in such a way maximizing also makes the agent achieve the goal. That is, the reward signal must not be used to impart any prior knowledge to the agent about how to achieve the goal. For example, in a chess game, rewards should be given when the agent wins the game, and not when accomplishing subgoals, such as taking enemy pieces. If the agent gets rewarded for achieving subgoals, it might find a way to maximize the rewards by only taking enemy pieces and without actually winning. All in all, the reward signal is our way of communicating to the agent what needs to be done, not how to do it.

# Chapter 3

## Markov Decision Processes

### 3.1 Markov Processes

In the reinforcement learning framework, the agent makes decisions as a function of a signal from the environment called the **state**. In this section, we discuss what is required of the state signal and what information it does or does not convey.

The state signal should include an immediate sensation, but it could include more than that, including some memory from past states. In fact, in typical applications, it usually is expected the state to inform the agent of more than just immediate sensations. For example, we could hear the word "yes", but we could be in totally different states depending on what was the question that came before and can no longer be heard. In contrast, the state signal should not be expected to inform the agent about all of the past experiences or all about the environment. Ideally, we want the state signal to summarize well past experiences, in such a way all the relevant information is retained.

To formalize this idea, for simplicity, suppose there are a finite number of states and rewards. Also, consider that the environment responds at time  $t + 1$  to an action taken at time  $t$ . We define the history sequence as  $h_t = \{S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t, S_t, A_t\}$ . Assume the state is a function of the history, that is,  $S_t = f(h_t)$ . If the state signal does not have the Markov property, the response of the environment depends on everything that happened earlier. Otherwise, the environment's response depends only on the state and actions at time  $t$ .

#### Definition 3.1: Markov Property

A state signal is said to have the Markov property if and only if

$$P\{R_{t+1} = r, S_{t+1} = s' \mid h_t\} = P\{R_{t+1} = r, S_{t+1} = s' \mid S_t, A_t\} \quad (3.1)$$

for all  $r, s'$  and  $h_t$ . If 3.1 is satisfied, then the environment also has the Markov property.

If every state of an environment is Markov, then we define a Markov Decision Process

(MDP) as follows:

### Definition 3.2: Markov Decision Process

A *Markov Decision Process* is a tuple  $\langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $S$  is a finite set of all valid states
- $\mathcal{A}$  is a finite set of all valid actions
- $P : S \times \mathcal{A} \rightarrow \mathcal{P}(S)$  is the transition probability function with  $P[S_{t+1} = s' \mid S_t = s, A_t = a]$  being the probability of transitioning into state  $s'$  starting in state  $s$  and taking action  $a$
- $\mathcal{R}$  is a reward function  
 $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$
- $\gamma$  is a discount factor such that  $\gamma \in [0, 1]$

Other authors also define an MDP to be a tuple  $\langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \rho_0 \rangle$ , with  $\rho_0$  being a distribution of starting states.

## 3.2 Policies

A big part of reinforcement learning is estimating how much reward the agent is expected to get by being in a specific state. In order to do that, we introduce the *value function* - a function of a state or a state-action pair that estimates the expected future rewards (or expected return), which tells us how good it is to be in a certain state or how good it is to take a specific action while in a specific state. Accordingly, value functions are defined with respect to different ways of acting and these ways of acting are dictated by a *policy*.

A policy is a mapping of actions to states, which can be either deterministic or stochastic. The first is a function  $\mu : S \rightarrow \mathcal{A}$  and the action at time  $t$  is

$$a_t = \mu(s_t)$$

and the latter is a probability distribution over  $a \in \mathcal{A}$  for each  $s \in S$  denoted by  $\pi$  and the action at time  $t$  is sampled from  $\pi$ :

$$a_t \sim \pi(\cdot \mid s_t)$$

Moreover, in this work, we will use *parameterized policies*, whose outputs are computable functions that depend on a set of parameters  $\theta$ , which can be adjusted using optimization algorithms. Parameterized policies are denoted by

$$\begin{aligned} a_t &= \mu_\theta(s_t) \\ a_t &= \pi_\theta(\cdot \mid s_t) \end{aligned}$$

Given an MDP  $\mathcal{M} = \langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and a policy  $\pi$ , we define the state-value function

**Definition 3.3: State-Value Function**

The *state-value function*  $v_\pi(s)$  of an MDP is the expected return starting from state  $s$  and following policy  $\pi$  afterwards

$$v_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \quad (3.2)$$

and the action-value function.

**Definition 3.4: Action-Value Function**

The *action-value function*  $q_\pi(s, a)$  of an MDP is the expected return starting from state  $s$ , taking action  $a$  and following policy  $\pi$  afterwards

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \quad (3.3)$$

The value function  $v_\pi$  can be estimated using *Monte Carlo Methods* by keeping the average of the returns that followed each state and then the average will eventually converge to the state's true value as the number of times that state is encountered approaches infinity. Similarly,  $q_\pi$  can be estimated by the same method by keeping the average of each state and each action taken in that state.

### 3.3 Optimality

Solving a reinforcement learning problem often means finding a policy that maximizes reward over the long run. Value functions define a partial ordering over policies:

$$\pi \geq \pi' \text{ if } v_\pi(s) \geq v_{\pi'}(s) \text{ for all } s \in \mathcal{S}$$

and then we are able to formulate what is an optimal policy.

**Theorem 3.1: Policy Optimality**

For any MDP, there exists an optimal policy  $\pi_*$  such that  $\pi_* \geq \pi, \forall \pi$ . All optimal policies achieve the optimal value function  $v_{\pi_*}$  and the optimal action-value function  $v_{\pi_*}$ .

**Definition 3.5: Optimal State-Value Function**

The optimal state-value function  $v_{\pi_*}$  or  $v_*(s)$  is the expected return starting from state  $s$  and always acting according to the optimal policy.

$$v_*(s) = \max_{\pi} v_\pi(s)$$

**Definition 3.6: Action-Value Function**

The optimal action-value function  $q_{\pi_*}$  or  $q_*(s)$  is the expected return starting from state  $s$ , taking an arbitrary action  $a$  and then always act according to the optimal policy.

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

Now, recall that the fundamental objective of reinforcement learning problems is to maximize rewards on the long run. One way to do that is using *Bellman Equations*, which enables us to calculate the optimal value functions defined above using recursive relationships and dynamic programming.

However, in this work, we'll focus on a particular class of methods called *Proximal Policy Optimization*, and for that, it isn't necessary to know how good a state/action is, but only how much better it is compared to others. We formalize this concept by defining the *advantage function*:

**Definition 3.7: Advantage Function**

The advantage function  $A^{\pi}(s, a)$  describes how much better it is to take a specific action  $a$  in state  $s$ , over randomly selecting an action according to  $\pi(\cdot | s)$  and following policy  $\pi$  forever after.

$$A_{\pi}(s, a) = q_{\pi}(s, a) - v_{\pi}(s)$$

### 3.4 Designing Rewards

One of the most important aspects when designing a RL experiment is defining reward functions. It's also one of the most challenging ones. For simpler problems, a positive reward +1 for achieving a goal and a negative reward -1 for not achieving the goal or losing, depending on the context. This is called a *sparse* reward, since rewards are only obtainable from a very specific and small set of states. In that way, while the agent is exploring the environment, it will not receive any feedbacks on how good or bad its actions were.

The most common solution is to just add more rewards via *reward shaping*. As described in Ng et al., 1999, reward shaping considers a *transformed MDP*  $M' = \langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}', \gamma \rangle$ , where  $R' = R + F$  and  $F: S \times \mathcal{A} \times S \rightarrow \mathbb{R}$  is a bounded real-valued function, called *reward shaping function*. While this approach may significantly improve performance and convergence, there is a possibility that we end up converging to a suboptimal policy and badly designed rewards can make the agent behave erratically or in an unintended way.

In Ng et al., 1999, the authors discuss a method called *potential-based* reward shaping (PBRS). This is the first method that has proven to leave the optimal policy unaltered (policy-invariance property). While PBRS was not used in our experiments, it is an important result in reward shaping in RL problems.



### Theorem 3.2: Potential-based shaping function

iven any  $S, \mathcal{A}, \gamma$  and shaping function  $F$ , we say  $F$  is a **potential-based** shaping function if there exists a real-valued function  $\Phi : S \rightarrow \mathbb{R}$  such that for all  $s \in S - \{s_0\}$ ,  $a \in \mathcal{A}$ ,  $s' \in S$ ,

$$F(s, a, s') = \gamma\Phi(s') - \Phi(s) \quad (3.4)$$

Then,  $F$  is a **sufficient** and **necessary** condition to guarantee optimal policy consistency when learning from MDP  $M' = \langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}', \gamma \rangle$ , rather than  $M = \langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ .

*Sufficiency* means that every optimal policy in  $M'$  is also an optimal policy in  $M$  and *necessity* means that if  $F$  is not a potential-based shaping function (there is no function  $\Phi$  satisfying 3.4), then there exist transition probability functions  $\mathcal{P}$  and reward functions  $\mathcal{R}$  such that no optimal policy in  $M'$  is an optimal policy in  $M$ .

The intuition behind PBRS is it avoids exploiting behaviors the agent may find in poorly designed rewards. Suppose a state  $s^*$  gives a reward of  $c \in \mathbb{R}$ , such that reaching this state is desirable and a state  $s'$  that gives no rewards. In an exploiting behavior, assuming the agent is already at  $s_t = s^*$ , it could transition to another state  $s_{t+1} = s'$  and go back to  $s_{t+2} = s^*$ , repeating this indefinitely to achieve more and more rewards. By using PBRS, this behavior would yield a negative reward:

$$F(s, a, s') = \gamma 0 - c = -c \quad (3.5)$$

That is, the agent would get punished for going back to state  $s'$ , because this would lead it to go to  $s^*$  again.



# Chapter 4

## Neural Networks

In this chapter, some concepts on *Neural Networks* (or NNs) will be discussed as well as a little bit of *Machine Learning*. As stated in the previous chapter, we'll be working with parameterized policies, and NNs will later be used as an approximator for said policies.

### 4.1 Machine Learning

Nowadays, as the Internet continues to expand, we have more and more information available and stored everyday. Due to that large volume of data, analyzing it and extracting meaningful information has become a task increasingly difficult for humans to perform. From that challenge, the concept of *Machine Learning* arises.

Machine learning can be defined as a set of techniques (or algorithms) that allows for a computer program to extract information, identify patterns and relationships in large amounts of data, in such a way a human cannot.

According to MITCHELL, 1997, "A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ".

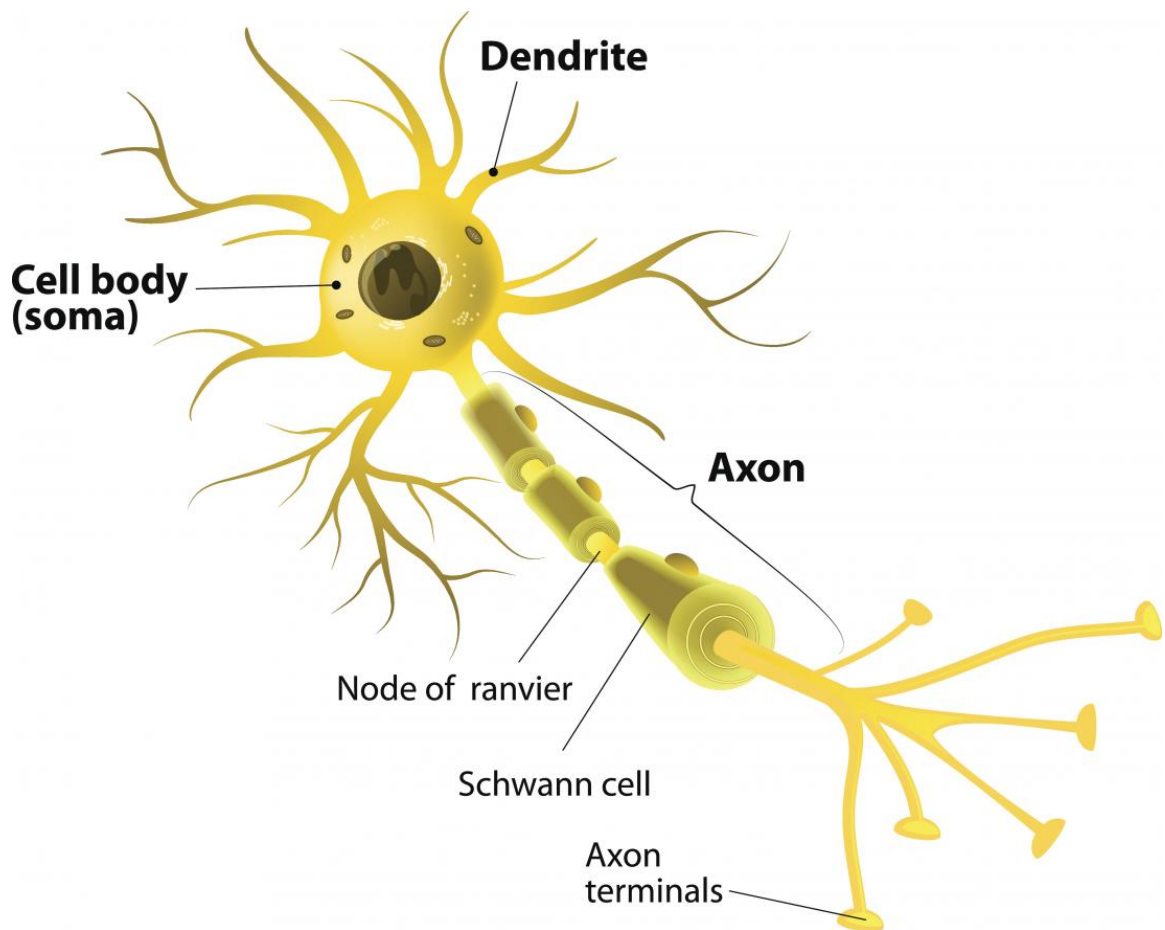
GOODFELLOW *et al.*, 2016 broke machine learning down to 3 main paradigms:

- In **Supervised Learning**, we have a dataset containing features, but each example is also associated with a *label* or *target*. For example, a dataset containing emails labeled as **spam** or **not spam**. A supervised learning algorithm can study this dataset and learn to classify whether an email is a spam or not.
- In **Unsupervised Learning**, we have a dataset containing many features, but no labels. Then, the goal is to learn useful properties of the structure of this dataset. In the context of deep learning, we usually want to learn the entire probability distribution that generated the dataset. Some other unsupervised learning algorithms perform other roles, like clustering, which consists of dividing the dataset into clusters of similar data.

- Then, there's **Reinforcement Learning**, where we don't have a fixed dataset. Instead, RL algorithms interact with the environment, forming a feedback loop between the learning system and what it has experienced.

## 4.2 Neural Networks

Neural Networks are a set of machine learning algorithms, whose structure is strongly related with the structure of the human brain.



**Figure 4.1:** Representation of a human neuron. Source: *MedicalNewsToday*

The main structural cells responsible for processing information are divided into 3 main parts:

- **Dendrites:** filaments responsible for receiving and transporting stimulus coming from the environment or other cells of the body.
- **Body cell:** processes the information gathered by the dendrites.
- **Axon:** transmits the nervous impulses generated by the cell body to the other cells through the axon terminals.

In summary, upon receiving stimuli, the cell body processes those signals and if the result exceeds a certain threshold value, the neuron fires a nerve impulse indicating it reacted to the input signals, which are further transmitted by other neurons to other neurons or cells.

The human brain processes information processes in order  $10^{-3}$  seconds, having a network of about 10 billion neurons densely connected, making it a huge, complex and efficient processing powerhouse, performing tasks like recognizing images and audios better than any machine. In an attempt to recreate a system that was able to operate like a human brain, the concept or *Artificial Neural Networks* was idealized.

**HAYKIN, 2009** cites the following useful properties and capabilities:

- **Nonlinearity:** an artificial neuron can be linear or nonlinear. Therefore, a neural network, made up of an interconnection of nonlinear neurons, is itself nonlinear. Nonlinearity is a highly important property, particularly if the underlying physical mechanism responsible for generation of the input signal (e.g., speech signal) is inherently nonlinear.
- **I/O Mapping:** neural networks perform very well in supervised learning tasks, where we have labeled data. Each piece of data is used to update the synaptic bias of the network such that the difference between the expected output and the actual output is minimized.
- **Adaptivity:** neural networks have the capability to adapt their synaptic weights to changes in the surrounding environment.
- **Evidential Response:** in particular, a neural network trained to operate in a specific environment can be easily retrained to deal with minor changes in the operating environmental conditions.
- **Contextual Information:** in the context of pattern classification, a neural network can be designed to provide information not only about which particular pattern to select, but also about the confidence in the decision made. This latter information may be used to reject ambiguous patterns, should they arise, and thereby improve the classification performance of the network.
- **Uniformity of Analysis and Design:** neural networks enjoy universality as information processors, in a sense that the same notation is used in all domains involving the application of neural networks, making it possible to share techniques and theories between models with different purposes.
- **Neurobiological Analogy:** the design of a neural network is motivated by analogy with the brain, which is a living proof that parallel processing is not only physically possible, but also fast and powerful.

## 4.3 Perceptron

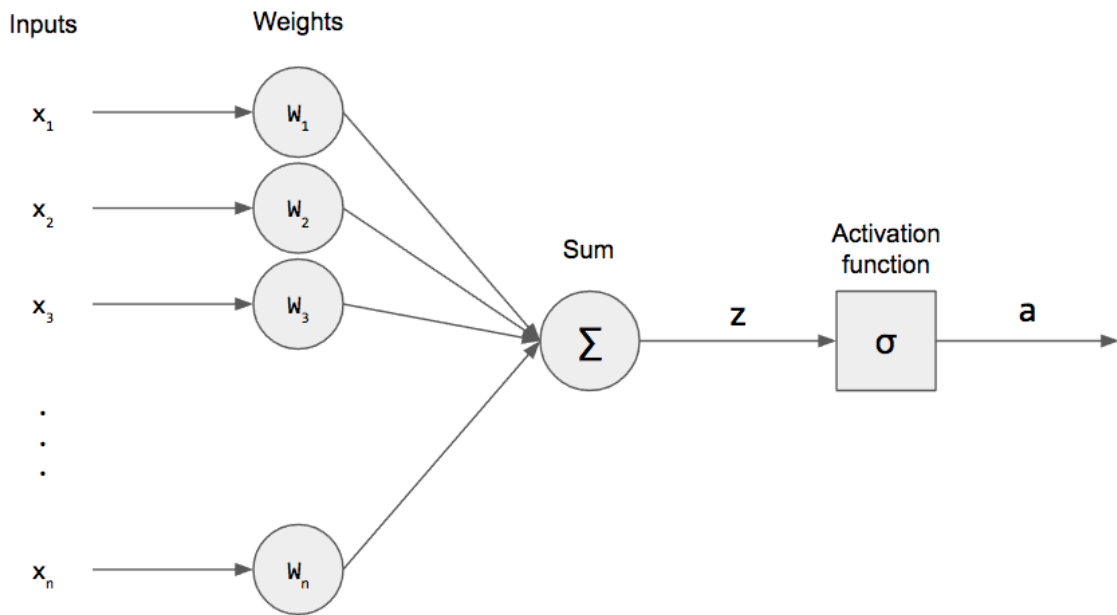
In order to create a model that would function similarly to a human brain, the psychologist and neurobiologist Frank Rosenblatt proposed the *perceptron* model. In his words, "a

probabilistic model for information storage and organization in the brain" [ROSENBLATT, 1958](#). In this artificial neuron, we have the following structures:

- An input vector  $x \in \mathbb{R}^n$
- The synaptic weights  $w \in \mathbb{R}^n$  responsible for weighing the values from the input vector. Large and positive values indicate higher relevance. Conversely, small and negative values indicate lower relevance.
- A bias  $b \in \mathbb{R}$ , which can be interpreted as how easily a neuron is activated.
- A linear combinator  $\Sigma$  weighing the input signals into a single scalar value:

$$z = \Sigma(x) = \sum_{i=1}^n w_i x_i + b$$

- An activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  mapping the weighted sum  $z$  to an output  $\phi(z)$ .



**Figure 4.2:** Representation of a artificial neural network. The bias has been omitted for brevity. Source: [Medium](#)

In general, given  $w \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ , the perceptron model can be defined as a function  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  such that

$$h(x \mid w, b) = \phi(w^T x + b)$$

From the perceptron, multiple models have been derived, but we are particularly interested in a specific *feed-forward* neural network, called *multilayer perceptron*.

## 4.4 Multilayer Perceptron

In a multilayer perceptron, the perceptrons (neurons) are stacked in multiple layers in such a way that every node on each layer is connected to all other nodes on the next layer, without any cycles, characterizing the *feed-forward* nature of the model.

The first layer is the input layer, and its units take the values of the input vector. The last layer is the output layer, and it has one unit for each value the network outputs. In the context of classification, it could have a single unit in the case of binary classification, or  $K$  units in the case of  $K$ -class classification. All the layers in between these are called hidden layers. It's called hidden because we don't know ahead of time what these units should compute, and this needs to be discovered during learning (GROSSE, 2021).

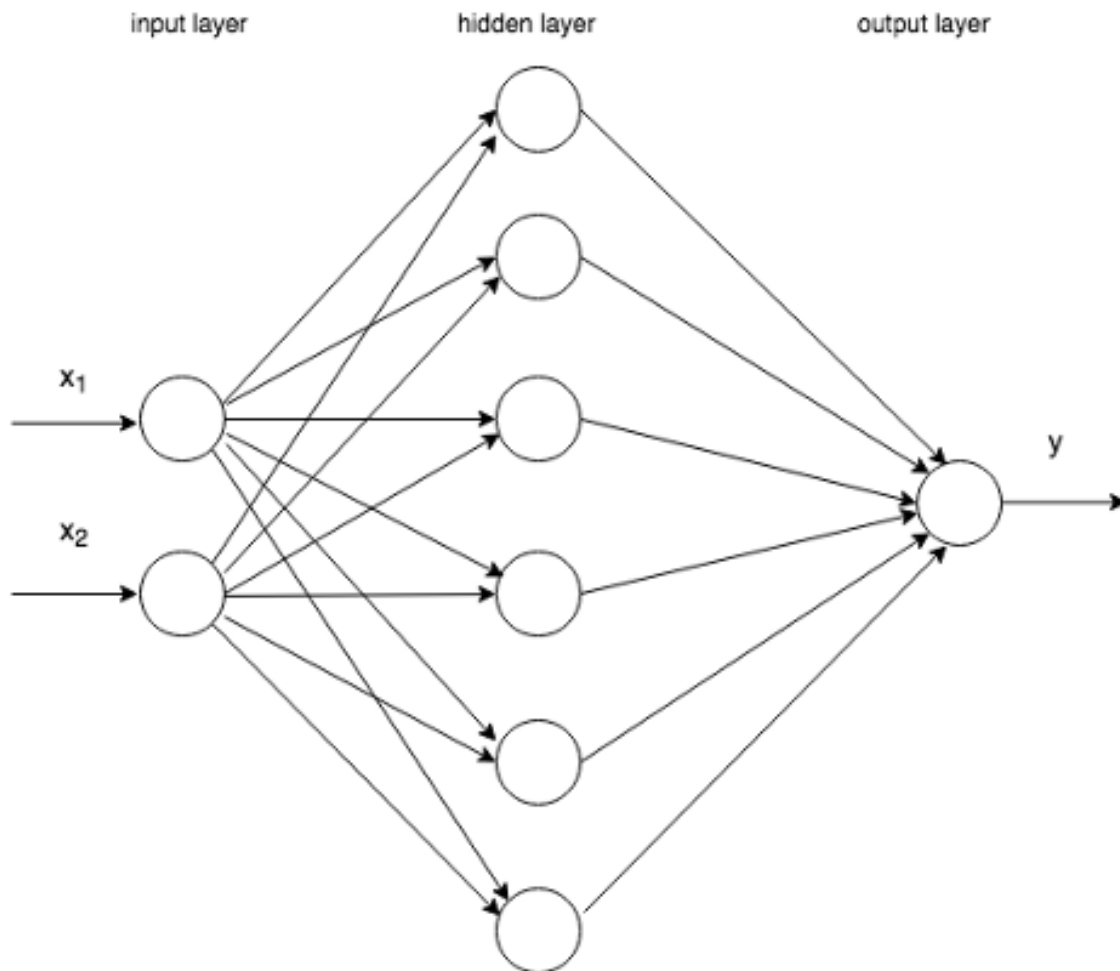


Figure 4.3: Representation of a multilayer perceptron. Source: *Medium*

Such as the perceptron, the MLP can also be defined as a function. Denote by  $h_i^{(l)}$  the  $i$ -th unit in the  $l$ -th hidden layer and by  $y$  the output unit. Note that each unit has its own bias, and there's a weight for every pair of units in two consecutive layers. Therefore, the

network's computation can be expressed as:

$$\begin{aligned} h_i^{(1)} &= \phi^{(1)} \left( \sum_j w_{ij}^{(1)} x_j + b_i^{(1)} \right) \\ h_i^{(2)} &= \phi^{(2)} \left( \sum_j w_{ij}^{(2)} h_j^{(1)} + b_i^{(2)} \right) \\ y_i &= \phi^{(3)} \left( \sum_j w_{ij}^{(3)} h_j^{(2)} + b_i^{(3)} \right) \end{aligned} \quad (4.1)$$

An important result for MLP models is the Universal Approximation Theorem (CYBENKO, 1989):

#### Theorem 4.1: Universal Approximation Theorem

Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous, bounded, non-constant and monotonically decreasing function. Let  $I_n = [0, 1]^n$  be the  $n$ -dimensional hypercube and  $C(I_n)$  the space of continuous functions in  $I_n$ . Then, given  $f \in C(I_n)$  independent of  $\sigma$  and  $\epsilon > 0$ , there exists  $N \in \mathbb{N}$ ,  $w_i \in \mathbb{R}^n$ ,  $b_i \in \mathbb{R}$  and  $\alpha_i \in \mathbb{R}$ , where  $i = 1, \dots, N$  such that

$$\begin{aligned} F(x) &= \sum_{i=1}^N \alpha_i \sigma(w_i x + b_i), \\ |F(x) - f(x)| &< \epsilon \end{aligned}$$

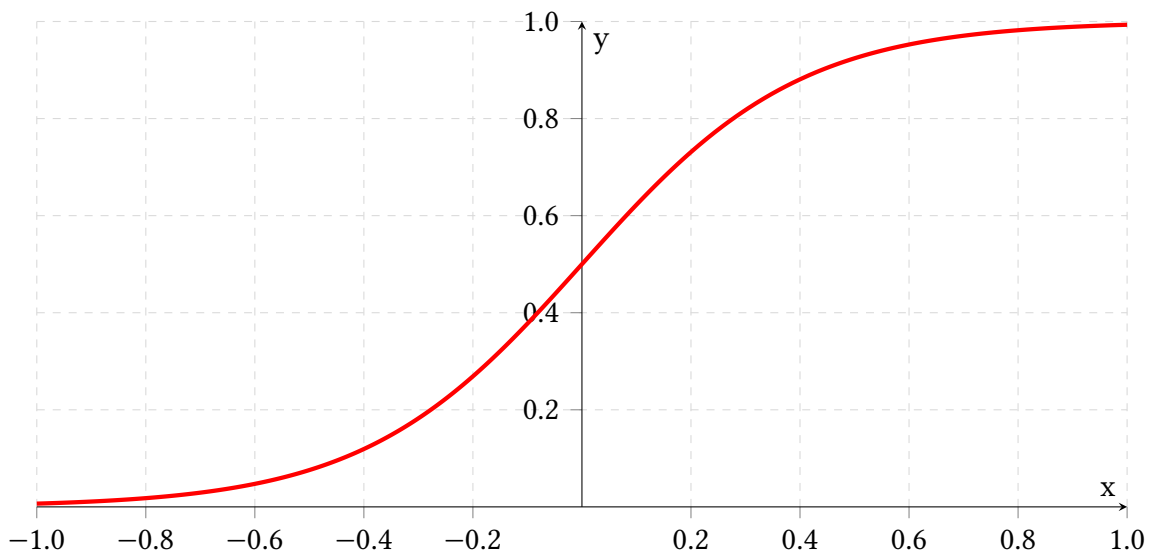
for all  $x \in I_n$ .

The theorem establishes that for any continuous function  $f$  on a compact subset of  $I_n$  can be approximated by a feed-forward neural network with only one hidden layer and finite number of units.

It's important to note that this doesn't imply that one neural network can accurately approximate any arbitrary continuous function under any circumstances. It is required that the neural network and its parameters (number of hidden units, number of learning iterations etc.) be adjusted for each unique function. Then, with the right parameters, it is possible to achieve any desired accuracy  $\epsilon$ . Cybenko proved the theorem specifically for the *sigmoid* activation function, defined as:

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$





**Figure 4.4:** Graph for the Sigmoid function. Source: Author.

But other activation functions have also been shown to satisfy the theorem, such as:

- Hyperbolic Tangent

$$\phi(z) = \frac{2}{1 + e^{-2z}} - 1 = 2\text{Sigmoid}(2z) - 1$$

- Softsign

$$\phi(z) = \frac{z}{1 + |z|}$$

Recently, some unbounded activation functions have also proven to be very effective as approximators, even while not satisfying the theorem:

- Rectified Linear Unit (ReLU)

$$\phi(z) = \max\{0, z\}$$

- Softplus

$$\phi(z) = \ln(e^z + 1)$$

- Exponential Linear Unit (ELU)

$$\phi(z) = \begin{cases} \alpha(e^z - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

for  $\alpha \in \mathbb{R}$ .

The theorem also has an extension for the case with multiple hidden layers and outputs, but will not be addressed here.



## Chapter 5

# Proximal Policy Optimization

Proximal policy optimization (PPO) was proposed by SCHULMAN, WOLSKI, *et al.*, 2017a as an alternative to already existing policy gradient (PG) methods, incorporating some concepts from trust region policy optimization (TRPO) methods, retaining some of its benefits while being significantly easier to implement.

### 5.1 Policy Gradient

For policy gradient, we consider parameterized policies, which can select actions without relying on a value function. The value function is still useful to learn the policy parameters, but it's not strictly necessary to select an action. This parameterization can be done in any way as long as the policy is differentiable with respect to its parameters.

Denoting by  $\theta \in \mathbb{R}^d$  the policy parameter vector, the probability of selecting action  $a$  at time  $t$  given that the environment is in state  $s$  with parameter  $\theta$  is

$$\pi_\theta(a \mid s, \theta) = P\{A_t = a \mid S_t = s, \theta_t = \theta\}$$

We also need to define a performance measure  $J(\theta)$  to quantify how good a policy is. We define such measure to be

$$J(\theta) = \sum_{s \in \mathcal{S}} d^{\pi_\theta}(s) \sum_{a \in \mathcal{A}} q_{\pi_\theta}(s, a) \pi_\theta(a \mid s) \quad (5.1)$$

where  $d^{\pi_\theta}(s) = \lim_{t \rightarrow \infty} P\{s_t = s \mid s_0, \pi\}$  is the stationary distribution of states under policy  $\pi$ , which is assumed to be independent of the starting state  $s_0$ . Policy gradient algorithms search for a local maximum in  $J$  using gradient ascent:

$$\Delta\theta = \alpha \nabla_\theta J(\theta)$$

where  $\nabla_{\theta}J(\theta)$  is the policy gradient defined as

$$\nabla_{\theta}J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

and  $\alpha$  is a step size parameter, commonly called *learning rate*.

The Policy Gradient Theorem provides a convenient way of expressing the gradient  $\nabla_{\theta}J(\theta)$  (adapted from [Richard S SUTTON et al., 1999](#)):

### Theorem 5.1: Policy Gradient Theorem

Given an MDP  $\langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and a parameterized policy  $\pi_{\theta}$ , the gradient of the expected return  $J(\theta)$  is given by

$$\nabla_{\theta}J(\theta) \propto \sum_{s \in S} d^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} q_{\pi_{\theta}}(s, a) \nabla_{\theta} \pi_{\theta}(a | s) \quad (5.2)$$

Furthermore, wW"log derivative trick", to rewrite the expression for the gradient:

$$\begin{aligned} \nabla_{\theta} \pi_{\theta}(a | s) &= \pi_{\theta}(a | s) \frac{\nabla_{\theta} \pi_{\theta}(a | s)}{\pi_{\theta}(a | s)} \\ &= \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s) \end{aligned}$$

Furthermore, using the identity above, we are able to express the gradient as an expectation:

$$\nabla_{\theta}J(\theta) \propto \sum_{s \in S} d^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} q_{\pi_{\theta}}(s, a) \nabla_{\theta} \pi_{\theta}(a | s) \quad (5.3)$$

$$\begin{aligned} &= \sum_{s \in S} d^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} q_{\pi_{\theta}}(s, a) \nabla_{\theta} \pi_{\theta}(a | s) \frac{\pi_{\theta}(a | s)}{\pi_{\theta}(a | s)} \\ &= \sum_{s \in S} d^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} q_{\pi_{\theta}}(s, a) \pi_{\theta}(a | s) \frac{\nabla_{\theta} \pi_{\theta}(a | s)}{\pi_{\theta}(a | s)} \\ &= \sum_{s \in S} d^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} q_{\pi_{\theta}}(s, a) \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s) \\ &= \mathbb{E}_{\pi} [q_{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a | s)] \end{aligned} \quad (5.4)$$

where  $\mathbb{E}_{\pi}$  is the expectation when  $s \sim d^{\pi_{\theta}}$  and  $a \sim \pi_{\theta}$ , i.e. both state and action distributions follow policy  $\pi_{\theta}$ . Expressing the gradient as an expectation means we can estimate it using a sample mean. We let the agent interact with the environment following a policy  $\pi_{\theta}$  and collect its *trajectory*  $\tau_i = \{s_0, a_0, \dots, s_{T+1}\}$  over  $N$  episodes, obtaining a set  $\mathcal{D} = \{\tau_i\}_{i=1, \dots, N}$  of trajectories. Then, the policy gradient is estimated as:

$$\hat{g} = \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^T q_{\pi_{\theta}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \quad (5.5)$$

That is, we compute the expression inside the expectation in each episode and take the sample mean as an estimator for the gradient, allowing us to take an update step.

Equation 5.4, also known as *grad-log-prob*, gives rise to an important result, which was used to derive a lot of other methods as an improvement over the *vanilla* policy gradient.

**Lemma 5.1: Expected Grad-Log-Prob (EGLP)**

Suppose that  $P_\theta$  is a parameterized probability distribution over a random variable  $x$ . Then:

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0 \quad (5.6)$$

*Proof.* First, recall that probability distributions are normalized:

$$\int_x P_\theta(x) dx = 1 \quad (5.7)$$

By taking the gradient on both sides, we get:

$$\nabla_\theta \int_x P_\theta(x) dx = \nabla_\theta 1 = 0 \quad (5.8)$$

Now, we can use the log derivative trick:

$$\begin{aligned} \nabla_\theta \int_x P_\theta(x) dx &= 0 \\ \int_x \nabla_\theta P_\theta(x) dx &= 0 \\ \int_x P_\theta(x) \nabla_\theta \log P_\theta(x) dx &= 0 \end{aligned} \quad (5.9)$$

□

As a consequence of the above lemma, SCHULMAN, MORITZ, *et al.*, 2015 proposed a more general form for policy gradients:

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \Phi_t \nabla_\theta \log \pi_\theta(a | s) \right] \quad (5.10)$$

where  $\Phi_t$  may be any of these functions:

- $\sum_{t=0}^{\infty} r_t$  : sum of total rewards.
- $\sum_{t'=t}^{\infty} r_{t'}$  : reward following action  $a_t$ .
- $q_\pi(s_t, a_t)$  : state-action value function.
- $A_\pi(s_t, a_t)$  : advantage function.

SCHULMAN, MORITZ, *et al.*, 2015 lists all the possible functions. For trust region methods and proximal policy optimization,  $\Phi_t$  is chosen to be the advantage function  $A_\pi(s_t, a_t)$ . The formulation of policy gradients with advantage functions is rather common, and the most known method to estimate it is *generalized advantage estimation* as described in SCHULMAN, MORITZ, *et al.*, 2015.

The policy gradient, while elegant, has shown to be problematic in practical problems:

**Sample inefficiency.** In order to run policy gradient, we need to sample from our policy millions and millions of times, since the estimation is done using Monte Carlo, averaging over a number of trial runs. Summing over all steps in a single trajectory could be very expensive computationally depending on the environment. It is also worth noting sample inefficiency is not a problem exclusive to policy gradient, it is an issue that has long plagued a lot of other RL algorithms.

**Slow convergence.** Sampling millions of trajectories is already inherently slow, and the high variance makes optimization very inefficient.

**High variance.** The high variance comes from the fact that, in RL, we are often dealing with very general problems. In our case, teaching a car to navigate through a parking lot. When we sample from an untrained policy, we are bound to observe highly variable behaviors, since we begin with a policy whose distribution of actions over states is effectively uniform. Of course, as the policy improves, the distribution is shaped to be unimodal or multimodal over some successful actions given a state. But in order to get there, we need the model to observe the outcomes of many different actions for each possible state. If we consider the action and state spaces to be continuous, the problem is even worse, since visiting every action-state pair possible is computationally intractable.

## 5.2 Trust Region Policy Optimization

TRPO updates policies by taking the largest step possible to improve performance, while satisfying a special constraint on how close the new and old policies are allowed to be.

Normal policy gradient keeps new and old policies close in parameter space, but even small differences in parameter space can have a large impact in performance, such that a single bad step can collapse the policy performance. Thus, it is dangerous to use large step sizes with vanilla policy gradients, which ultimately makes the method very sample inefficient. TRPO not only avoids this kind of collapse, but also tends to quickly and monotonically improve performance.

The way TRPO achieves this is by guaranteeing that the policy doesn't change too much in comparison to the old one using Kullback-Leibler divergence. KL divergence is a statistical measure of how different a probability distribution is from another.

The objective function in TRPO is

$$J(\theta) = \mathbb{E}_{a \sim \pi_{\theta_{\text{old}}}} \left[ \frac{\pi_{\theta}(a | s)}{\pi_{\theta_{\text{old}}}(a | s)} \hat{A}_{\pi_{\theta_{\text{old}}}} \right] \quad (5.11)$$

and the goal is to maximize it subject to the *trust region constraint*, which enforces the distance between old and new policies, measured by KL-divergence, to be small enough:

$$\mathbb{E}[\text{KL}(\pi_{\theta_{\text{old}}}(a | s), \pi_{\theta}(a | s))] \leq \delta \quad (5.12)$$

That way, the old and new policies would not diverge too much when this hard constraint is met. Not only that, but TRPO also guarantees monotonic improvement over each iteration.

The full details of the derivation of this method have been omitted, but can be found in SCHULMAN, LEVINE, *et al.*, 2015.

The same way TRPO emerged as an improvement over vanilla policy gradient, PPO emerges as an improvement over TRPO. In short, some of the disadvantages of TRPO is that it is computationally expensive, still sample inefficient and its derivation is complex.

## 5.3 Proximal Policy Optimization

Proximal policy optimization was proposed in SCHULMAN, WOLSKI, *et al.*, 2017b and the objective function is similar to TRPO, and is called *clipped surrogate objective function*:

$$L^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip} \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (5.13)$$

where  $r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$ . This is the same ratio as seen in TRPO and is one way to measure the divergence between the old and current policy. Note that if  $r_t(\theta) > 1$ , then, at time  $t$ , taking action  $a$  in state  $s$  is more likely in the current policy than the old policy. If  $0 < r_t(\theta) < 1$ , then the action is less likely in the current policy than the old one. According to SCHULMAN, WOLSKI, *et al.*, 2017b, the use of this ratio has been originally proposed by KAKADE and LANGFORD, 2002.

As for the second term inside the min function, instead of using KL divergence to limit the distance between the two policies, PPO instead uses the clip operator, keeping that distance between the interval  $[1 - \epsilon, 1 + \epsilon]$ . The clip operator is defined as

$$\text{clip} \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) = \begin{cases} 1 - \epsilon & \text{if } r_t(\theta) < 1 - \epsilon \\ 1 + \epsilon & \text{if } r_t(\theta) > 1 + \epsilon \\ r_t(\theta) & \text{otherwise} \end{cases} \quad (5.14)$$

where  $\epsilon$  is a hyperparameter. In SCHULMAN, WOLSKI, *et al.*, 2017b and in a lot of practical applications,  $\epsilon = 0.2$  is used. Clipping (or restricting) the range of values that the probability ratio  $r_t(\theta)$  can take is supposed to remove the incentive for pushing the probability ratio outside the interval enforced by the clipping operator. In other words, due to the clipping

operation, the probability ratio  $r_t(\theta)$  is supposed to remain within the interval  $[1 - \epsilon, 1 + \epsilon]$  even after multiple updates. This way, we avoid destructively large weight updates, but in a different fashion than TRPO. Finally, we take the minimum of the clipped and the non-clipped objective, such that the final objective is a lower bound of the clipped objective.

In cases where clipping does not apply, that is, when the value of  $r_t(\theta)$  lies within the interval  $[1 - \epsilon, 1 + \epsilon]$ , neither the minimum operator or the clip operator impact the computation of the gradient.

Now, consider the cases where clipping applies. When this happens, the behavior of  $L^{\text{CLIP}}$  depends on the sign of  $\hat{A}_t$ .

If  $r_t(\theta) < 1 - \epsilon$  and  $\hat{A}_t > 0$ , the minimum operator will always select the value of  $r_t(\theta)\hat{A}_t$  instead of the value computed by the clip operator, since

$$\begin{aligned} r_t(\theta) &< 1 - \epsilon < \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) = 1 - \epsilon \\ r_t(\theta)\hat{A}_t &< (1 - \epsilon)\hat{A}_t < \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t = (1 - \epsilon)\hat{A}_t \end{aligned}$$

The probability function being less than  $1 - \epsilon$  means that the probability of choosing some action  $a_t$  in state  $s_t$  has decreased during the previous weight updates, but it was better than expected, as indicated by  $\hat{A}_t > 0$ . Thus, the gradient will make so that action  $a_t$  becomes more likely in state  $s_t$  again.

If  $r_t(\theta) < 1 - \epsilon$  and  $\hat{A}_t < 0$ , then the inequality sign flips

$$\begin{aligned} r_t(\theta) &< 1 - \epsilon < \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) = 1 - \epsilon \\ r_t(\theta)\hat{A}_t &> (1 - \epsilon)\hat{A}_t > \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t = (1 - \epsilon)\hat{A}_t \end{aligned}$$

and the minimum operator selects the clipped objective, which evaluates to  $(1 - \epsilon)\hat{A}_t < 0$ . Since this is a constant (because it's not a function of  $\theta$ ), the gradient will evaluate to zero. Note that if this wasn't the case, the gradient would make the probability of selecting action  $a_t$  in state  $s_t$  even more unlikely, increasing the divergence between the two policies and ultimately collapsing performance, because the updates would be large, and this is exactly what PPO tries to avoid.

The cases when  $r_t(\theta) > 1 + \epsilon$  are analogous. All the non-trivial cases are summarized in the table below:

$r_t(\theta)$	$A_t$	min return value	Objective clipped?	Objective sign	Gradient
$r_t(\theta) \in [1 - \epsilon, 1 + \epsilon]$	+	$r_t(\theta)A_t$	no	+	non-zero
$r_t(\theta) \in [1 - \epsilon, 1 + \epsilon]$	-	$r_t(\theta)A_t$	no	-	non-zero
$r_t(\theta) < 1 - \epsilon$	+	$r_t(\theta)A_t$	no	+	non-zero
$r_t(\theta) < 1 - \epsilon$	-	$(1 - \epsilon)A_t$	yes	-	<b>0</b>
$r_t(\theta) > 1 + \epsilon$	+	$(1 + \epsilon)A_t$	yes	+	<b>0</b>
$r_t(\theta) > 1 + \epsilon$	-	$r_t(\theta)A_t$	no	-	non-zero



# Chapter 6

## Experiments

The main goal of this work was to use proximal policy optimization to simulate a parking situation in a virtual environment.

The tool used to create the simulation was Unity (version 2021.3.7f1), developed by Unity Technologies, commonly used as a game engine. All the assets used are available for free in the Unity Store.

Unity provides a open-source *toolkit* called ML-Agents, which enables developers to create and train AI agents inside the platform. ML-Agents also provides its own implementation of proximal policy optimization, as described in [SCHULMAN, WOLSKI, et al., 2017b](#).

The first experiment is the simplest, where we keep every spawn fixed across all episodes. In the following experiments, we randomize the spawns of the agent, parking spot and obstacles (parked cars), until the environment is completely random.

### 6.1 Hyperparameters and Evaluation Metrics

In the context of machine learning, a *hyperparameter* is a parameter whose value is used to control the learning process. The main difference of a hyperparameter and a parameter is that the values of hyperparameters are manually set, while parameters (for example, node weights in a neural network) are derived during training.

For the neural network, which will be used as an approximator for the policy, the hyperparameters to be set are the **number of hidden layers** and the **number of units per layer**, as seen in section 4.2.

For PPO, we need to set the **batch size**, **buffer size**, **divergence limit**  $\epsilon$ , **exponential weight discount**  $\lambda$ , **learning rate**  $\alpha$ , **number of epochs**, **time horizon** and **max step**.

**Batch size.** Determines how many experiences (trajectories) are used for each gradient descent update.

**Buffer size.** Determines how many experiences (trajectories) are collected before a gradient descent update is performed on them all.

**Epochs.** It is the number of passes through the experience buffer during gradient descent.

We update the model every time we get to a determined number of timesteps, determined by the buffer size. During this update, the buffer is divided into  $n$ -sized batches, determined by batch size, and a gradient descent update is performed on each of these batches one at a time. The number of times this process is repeated is determined by the number of epochs.

**Time horizon.** Determines how many steps of experience to collect before adding it to the experience buffer. This number should be high enough to allow the agent to explore the environment and capture all the important behaviors within a sequence of actions.

**Max step.** Determines how many steps of the simulation are to be run during the training process.

**Divergence limit.** As described in section 5.3, this is a threshold on the divergence between the old and new policies. Setting a lower value results in more stable updates, but may significantly slow training process.

**Exponential weight discount.** Needed for generalized advantage estimation, the method used to estimate the advantage function used in PPO. For details, see [SCHULMAN, MORITZ, \*et al.\*, 2015](#).

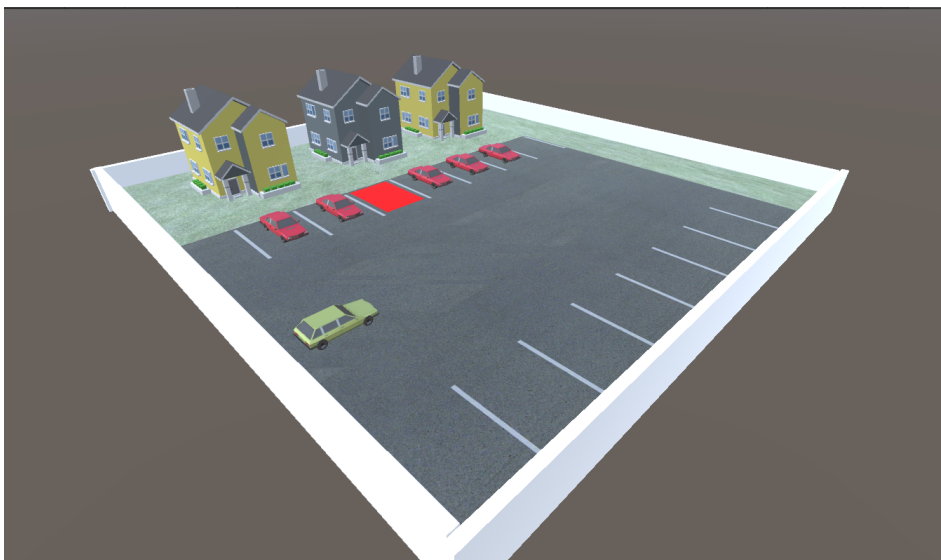
**Learning rate.** Determines the step size of a gradient descent update. Larger values may collapse performance and smaller values may increase training time.

As for evaluation metrics, we are usually most concerned with the cumulative reward, since the ultimate goal is to maximize it. But other metrics can be useful to either verify how the policy is changing over time or diagnose training issues. For our experiments, we'll also be displaying the episode length and entropy, a measure of how random decisions are. If entropy is high, certainty is low and vice versa.

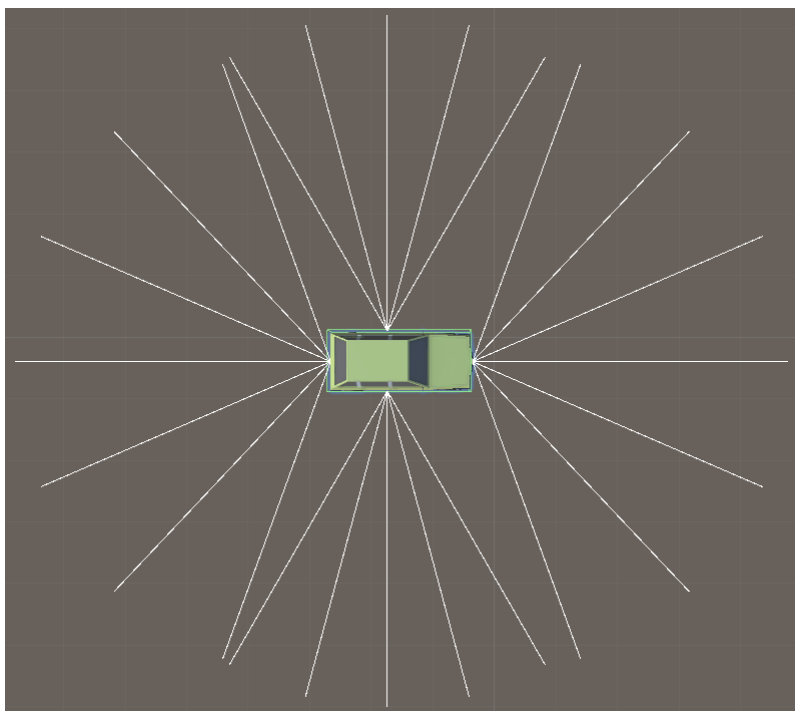
## 6.2 Experiment 1: Fixed Positions

We design a fairly simple parking scenario with a few cars already parked serving as obstacles. The environment is a one floor parking scenario with a total of 16 parking spots. For this first experiment, the designated parking spot is fixed, as well as the other cars' positions. The car is considered to be parked when it stays within 0.6 units from the spot for at least 0.5 seconds. A unit in Unity is equivalent to 1 meter and the distance is computed from the center of the car to the center of the parking spot.

The agent is a low-polygon 3D car model equipped with a total of 24 depth sensors which can detect objects up to 5 units away.



**Figure 6.1:** *The parking lot environment in Unity Engine*



**Figure 6.2:** *The car (agent) and its sensors*

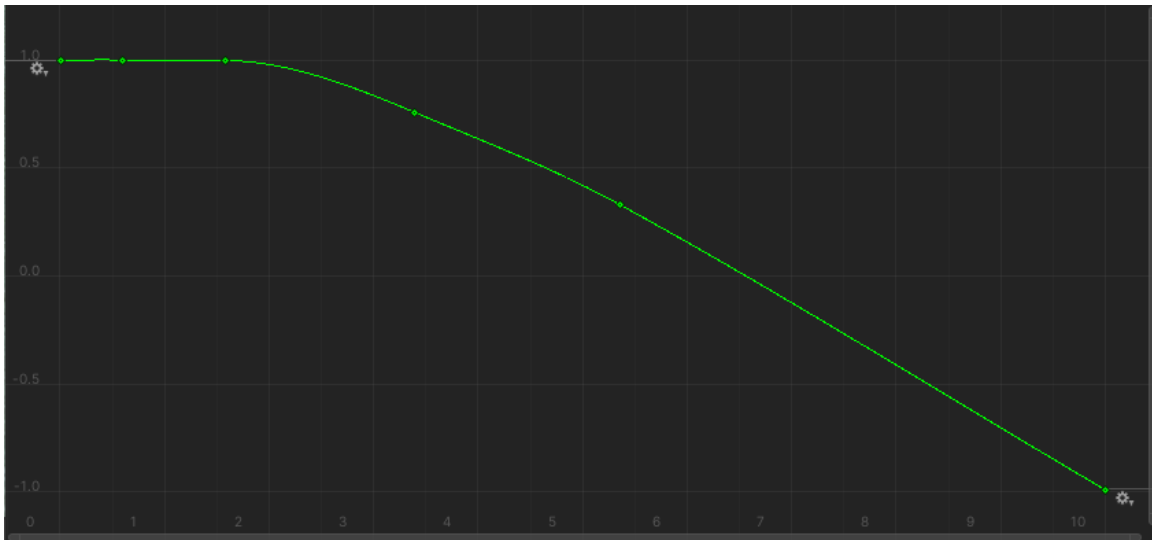
The reward functions are described in the table below:

Reward function	Value
Timeout	−1000
Parking successfully	1000
Parking successfully and sufficiently aligned	5000
Time	−0.1/timestep
Stopping	−0.0002/timestep
Collision	−1000
Stay in collision	−0.5/timestep
Goal heading	−1 to 1/timestep, see below for details
Goal distance	−1 to 1/timestep, see below for details

**Table 6.1:** *Reward functions for the agent.*

During the first few experimental runs, we tried multiple ways of encouraging the agent to get closer to the goal, but, of course, without actually telling the exact coordinates. Rewarding it for getting closer compared to the previous time step was our first successful attempt at teaching the agent to park at the designated spot, but the success rate was still rather low. The agent would often get stuck in a loop going back and forth and the episode would eventually end due to timeout.

With the intent of providing a better incentive for getting closer, we designed a custom curve to define the rewards the agent would get based on its distance from the parking spot.



**Figure 6.3:** *The reward function for getting closer to the goal*

The  $x$ -axis ranges from 0 to 10 and represents the distance between the agent and the goal. The  $y$ -axis ranges from −1 to 1 and increases as the agent gets closer to the spot. From the very first attempt, these reward functions showed very good results, with success rates above 50%. However, the agent would often park completely disaligned.

To seek further improvement, we also rewarded the agent at each time step based on whether it was heading towards the goal or not. The reward is the dot product between the vector pointing forward from the car and the vector pointing forward from the parking spot, such that the reward is 1 when the agent and the parking spot are perfectly aligned and  $-1$  when the agent is headed the complete opposite way. Note that this measure is invariant of distance. Adding this extra reward function not only increased the success rate to almost 100%, but also significantly decreased training time. Unfortunately, it hadn't solved the original issue.

As an second attempt, we added an extra criterion to define whether the agent is parked or not. If the dot product between the two vectors mentioned above is  $\geq 0.9$  at the moment of parking, the reward of 1000 is given and the episode ends. If the dot product is  $\geq 0.97$  the reward is increased to 5000. In angles, since both vectors are unitary, the dot product is simply the cosine between the two vectors, that is, the angle must be at most  $45^\circ$  to consider the agent parked and at most  $\approx 25^\circ$  to receive the bonus reward. In terms of training time and success rate, nothing has changed, but the agent would always park correctly and get the bonus once trained.

Lastly, we noted that the policy the model converged to was not always the same. One of these policies is what we actually intended to achieve - parking as aligned as possible and in the least time possible. The second is what we believe to be the optimal policy, since it yielded the most reward compared to the other - the agent would stop very close to the spot, in such a way to receive the maximum reward possible from both the alignment and distance criteria. The workaround was to increase the penalty at each time step, so that exploiting the rewards for distance and alignment would never be more beneficial than parking. After that, the model would always converge to the same policy, guaranteeing the result was reproducible.

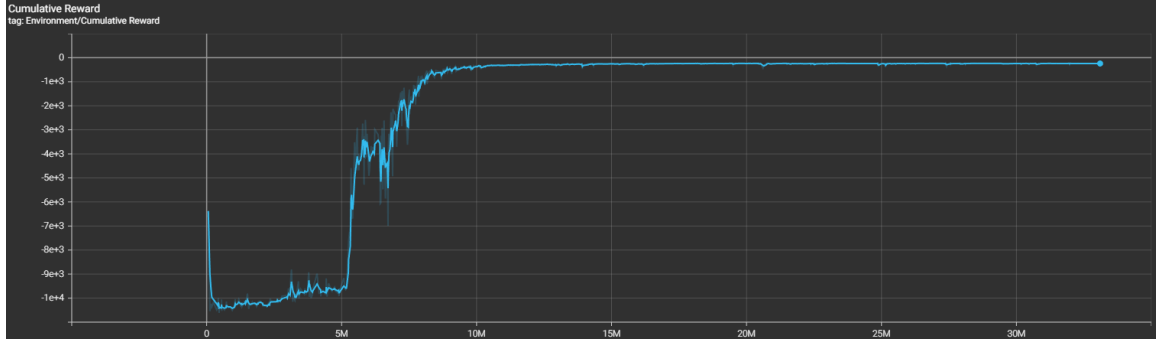
The value of hyperparameters used and the results are presented below.

Hyperparameter	Value
Divergence limit $\epsilon$	0.2
Exponential discount $\lambda$	0.95
Epochs	3
Discount factor $\gamma$	0.99
Batch size	2048
Buffer size	20480
Learning rate $\alpha$	$3 \cdot 10^{-4}$
Number of layers	2
Number of hiddens unit	256
Time horizon	512
Max step	$5 \cdot 10^7$

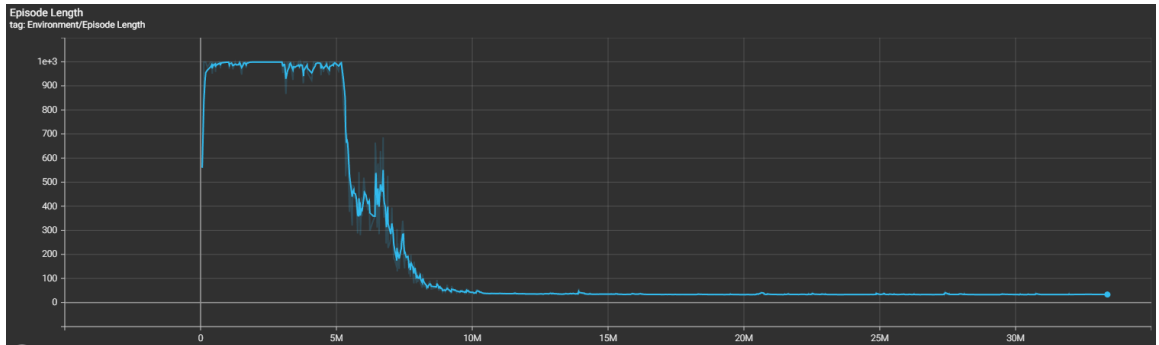
**Table 6.2:** *Hyperparameter configuration.*

The values used are the suggested values for training continuous action spaces models by Unity ML-Agents official documentation. The only hyperparameters tweaked were the

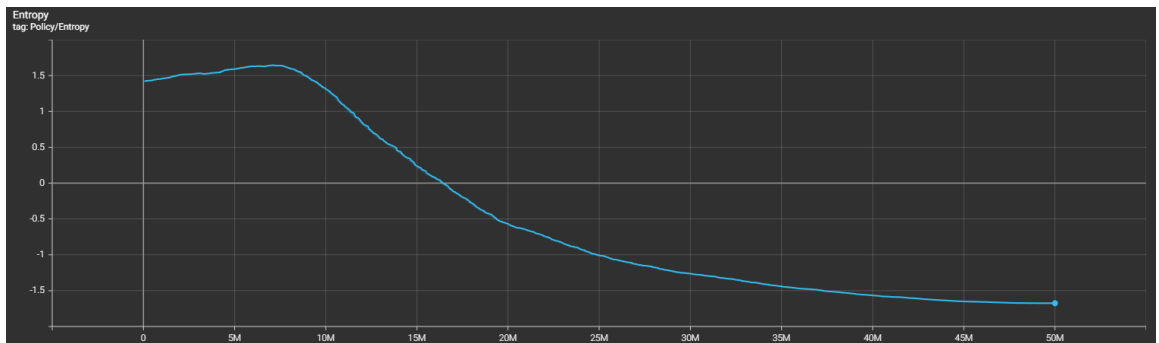
time horizon and the max step, to guarantee the agent would have enough time to explore the environment and to make sure we would have enough training steps to converge to a policy.



**Figure 6.4:** Cumulative rewards obtained by the agent during training.



**Figure 6.5:** Episode length during training.

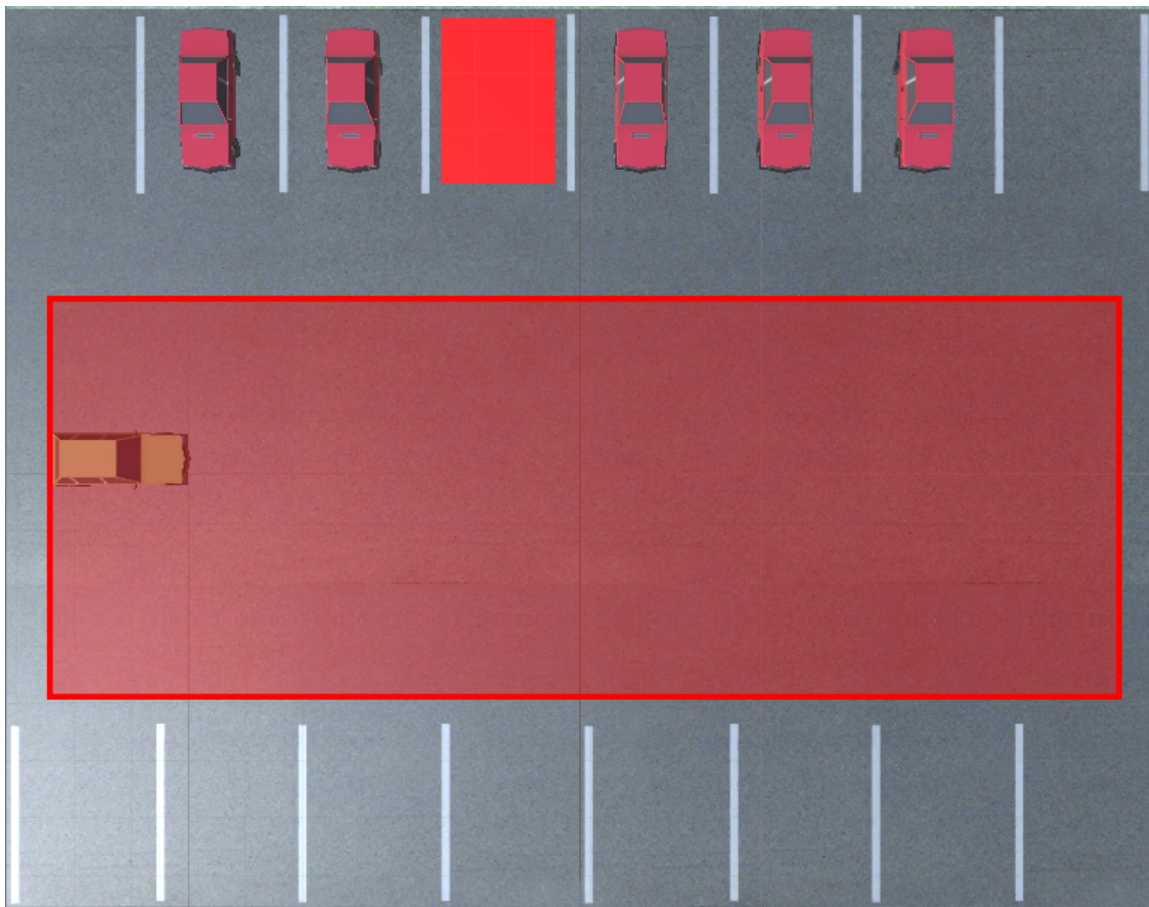


**Figure 6.6:** Entropy evolution during training.

The behavior is completely random when training starts, but once it parks correctly for the first time, the agent is quick to learn what is to be achieved, and eventually learns how to do it in the least time possible. To summarize, the main challenges of this experiment was to provide enough incentive for the agent to reach the goal and working around eventual exploits the agent would find in our reward design that would make it not accomplish what was intended.

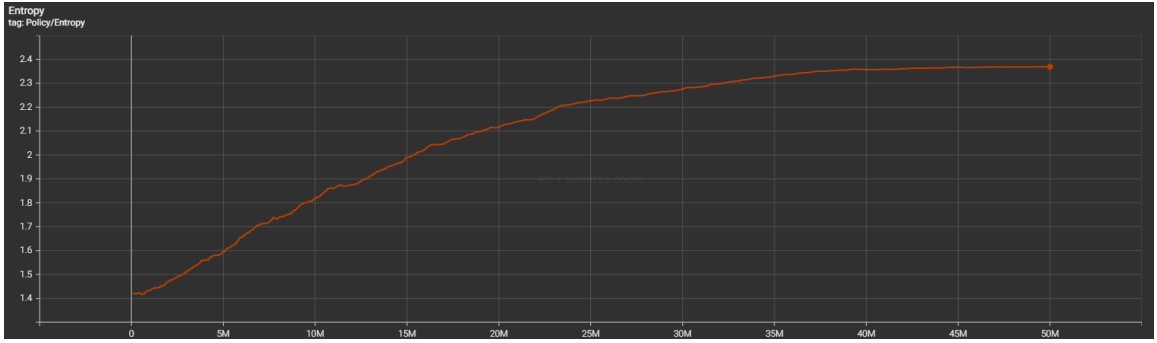
## 6.3 Experiment 2: Randomized Car Position and Fixed Parking Spot

For this experiment, we added some randomness to the environment by randomizing the starting position of the agent. During the first iteration of this experiment, we selected an area within the parking lot and at each episode, the agent would start in a random position within that area.



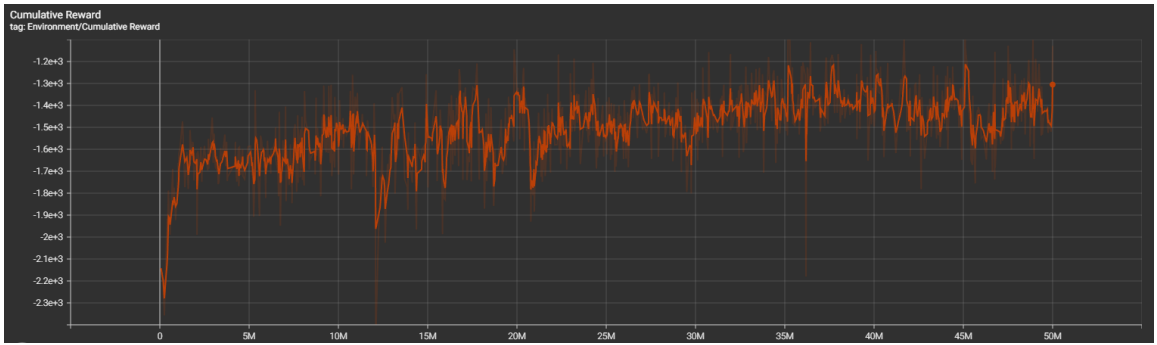
**Figure 6.7:** Area for the agent's starting position.

All hyperparameters and rewards were kept unchanged. Ultimately, in this setup, we couldn't reach the desired behavior, with the agent's actions being unintelligible even after  $5 \cdot 10^7$  iterations, as shown by the policy entropy



**Figure 6.8:** *Entropy evolution during training.*

and the cumulative rewards having very high variance, compared to the previous experiments.



**Figure 6.9:** *Cumulative rewards obtained by the agent during training.*

We gradually reduced the area in which the agent could spawn in an attempt to reduce randomness and achieve convergence, but to no avail.

Ultimately, we decided to work with a finite amount of possible initial positions, until we got to only 2 possible positions and still did not achieve convergence. While the model showed an increase in cumulative rewards, it still had high variance.



# Chapter 7

## Conclusions

In this work, we presented an overview of how reinforcement learning works, with a particular focus on policy gradient methods and its derivation, showing some of its pitfalls and how this motivated the creation of TRPO, to address stability issues found in vanilla PG. Even being more computationally expensive, TRPO is still substantially more efficient when working with high-dimensional problems, effectively balancing out with the computational cost drawback.

Then, PPO was proposed to address the high computational cost of TRPO while retaining similar performance levels and being more sample efficient. While PPO is often considered as one of the most important recent breakthroughs in the field, a recent study (ENGSTROM *et al.*, 2020) has rised some controversy by contradicting the belief that replacing the KL divergence found in TRPO by the clipping operator is the key innovation in PPO. The author also concludes that it's still poorly understood how each component of a RL algorithm impacts agent training performance.

Controversies apart, proximal policy optimization showed good results in our experiments, however, contrary to our initial belief when first designing the experiments, properly designing rewards proved to be one of the most important aspects of a RL experiment. It was easy to see how the policy completed changed as we tried different reward functions for the agent. This motivated further study in design rewards, which ultimately ended up in its own section (3.4).

As for future extensions of this work, the most important one is to refine the rewards design, such as rewriting some of them as a potential-based reward function to guarantee policy invariance, as well as come up with new ones.

Increasing complexity is also a possibility, adding more obstacles, including moving ones, such as pedestrians, or adding more than one possible parking spot and observe if the agent would be able to consistently go for the closest one.

It could also be interesting to include other parking situations we found on our daily lives, such as parallel parking or reverse parking.



# References

- [CYBENKO 1989] George CYBENKO. “Approximation by superpositions of a sigmoidal function.” In: *Math. Control. Signals Syst.* 2.4 (1989), pp. 303–314. URL: <http://dblp.uni-trier.de/db/journals/mcss/mcss2.html#Cybenko89> (cit. on p. 18).
- [ENGSTROM *et al.* 2020] Logan ENGSTROM *et al.* *Implementation Matters in Deep Policy Gradients: A Case Study on PPO and TRPO*. 2020. DOI: [10.48550/ARXIV.2005.12729](https://arxiv.org/abs/2005.12729). URL: <https://arxiv.org/abs/2005.12729> (cit. on p. 35).
- [GOODFELLOW *et al.* 2016] Ian GOODFELLOW, Yoshua BENGIO, and Aaron COURVILLE. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016 (cit. on p. 13).
- [GROSSE 2021] Roger GROSSE. *Lecture 5: Multilayer Perceptrons*. URL: [https://www.cs.toronto.edu/~rgrosse/courses/csc311\\_f21/](https://www.cs.toronto.edu/~rgrosse/courses/csc311_f21/). 2021 (cit. on p. 17).
- [HAYKIN 2009] Simon S. HAYKIN. *Neural networks and learning machines*. Third. Pearson Education, 2009 (cit. on p. 15).
- [KAKADE and LANGFORD 2002] Sham KAKADE and John LANGFORD. “Approximately optimal approximate reinforcement learning”. In: *IN PROC. 19TH INTERNATIONAL CONFERENCE ON MACHINE LEARNING*. 2002, pp. 267–274 (cit. on p. 25).
- [MITCHELL 1997] Tom M. MITCHELL. *Machine Learning*. McGraw-Hill, 1997 (cit. on p. 13).
- [NG *et al.* 1999] Andrew Y. NG, Daishi HARADA, and Stuart J. RUSSELL. “Policy invariance under reward transformations: theory and application to reward shaping”. In: *Proceedings of the Sixteenth International Conference on Machine Learning*. ICML ’99. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1999, pp. 278–287. ISBN: 1558606122 (cit. on p. 10).
- [ROSENBLATT 1958] Frank ROSENBLATT. “The perceptron: A probabilistic model for information storage and organization in the brain”. In: *Psychological Review* 65 (1958) (cit. on p. 16).

- [SCHULMAN, LEVINE, *et al.* 2015] John SCHULMAN, Sergey LEVINE, Philipp MORITZ, Michael I. JORDAN, and Pieter ABBEEL. “Trust region policy optimization”. In: *CoRR* abs/1502.05477 (2015). arXiv: [1502.05477](https://arxiv.org/abs/1502.05477). URL: <http://arxiv.org/abs/1502.05477> (cit. on p. 25).
- [SCHULMAN, MORITZ, *et al.* 2015] John SCHULMAN, Philipp MORITZ, Sergey LEVINE, Michael JORDAN, and Pieter ABBEEL. *High-Dimensional Continuous Control Using Generalized Advantage Estimation*. 2015. DOI: [10.48550/ARXIV.1506.02438](https://doi.org/10.48550/ARXIV.1506.02438). URL: <https://arxiv.org/abs/1506.02438> (cit. on pp. 23, 24, 28).
- [SCHULMAN, WOLSKI, *et al.* 2017a] John SCHULMAN, Filip WOLSKI, Prafulla DHARIWAL, Alec RADFORD, and Oleg KLIMOV. *Proximal Policy Optimization Algorithms*. 2017. DOI: [10.48550/ARXIV.1707.06347](https://doi.org/10.48550/ARXIV.1707.06347). URL: <https://arxiv.org/abs/1707.06347> (cit. on p. 21).
- [SCHULMAN, WOLSKI, *et al.* 2017b] John SCHULMAN, Filip WOLSKI, Prafulla DHARIWAL, Alec RADFORD, and Oleg KLIMOV. *Proximal Policy Optimization Algorithms*. 2017. arXiv: [1707.06347](https://arxiv.org/abs/1707.06347) [cs.LG] (cit. on pp. 25, 27).
- [Richard S SUTTON *et al.* 1999] Richard S SUTTON, David MCALLESTER, Satinder SINGH, and Yishay MANSOUR. “Policy gradient methods for reinforcement learning with function approximation”. In: *Advances in Neural Information Processing Systems*. Ed. by S. SOLLA, T. LEEN, and K. MÜLLER. Vol. 12. MIT Press, 1999. URL: <https://proceedings.neurips.cc/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf> (cit. on p. 22).
- [Richard S. SUTTON and BARTO 2015] Richard S. SUTTON and Andrew G. BARTO. *Reinforcement Learning: An Introduction*. The MIT Press, 2015 (cit. on pp. 1, 5, 6).