



POLITECNICO DI TORINO

Department of Control and Computer Engineering

Master of Science in Computer Engineering (Software Career)

Master Thesis

# Deep Reinforcement Learning algorithms for autonomous systems

Design and implementation of a control system for autonomous  
driving task of a small robot, exploiting state-of-the-art Model-Free  
Deep Reinforcement Learning algorithms

## **Supervisors**

prof. Pietro MICHIARDI

prof. Elena BARALIS

## **Candidate**

Piero MACALUSO

matricola: s252894

ACADEMIC YEAR 2019-2020

This work is subject to the Creative Commons Licence

# Contents

<b>Contents</b>	III
<b>1 Introduction</b>	1
<b>2 Reinforcement Learning</b>	2
2.1 Fundamentals of reinforcement learning . . . . .	2
2.1.1 The reinforcement learning problem . . . . .	3
2.1.2 Bellman equations . . . . .	6
2.1.3 Approaches to Reinforcement Learning . . . . .	7
2.1.4 Dynamic programming . . . . .	9
2.1.5 Model-free approach . . . . .	10
2.1.6 Model-based approach . . . . .	14
2.2 Deep reinforcement learning . . . . .	14
2.2.1 Fundamentals of Deep Learning . . . . .	15
2.2.2 Value-based methods . . . . .	22
2.2.3 Policy gradient methods . . . . .	25
2.2.4 Deep Deterministic Policy Gradient (DDPG) . . . . .	27
2.2.5 Soft Actor-Critic (SAC) . . . . .	30
2.3 Related Work . . . . .	34
<b>3 Tools and Frameworks</b>	38
3.1 OpenAI Gym . . . . .	38
3.1.1 Environments . . . . .	39
3.1.2 Observations . . . . .	41
3.2 Anki Cozmo . . . . .	42
3.2.1 Cozmo Architecture . . . . .	43
3.2.2 Why Cozmo? . . . . .	46
3.3 PyTorch . . . . .	49
3.3.1 TensorboardX . . . . .	51
3.3.2 PyTorch vs. TensorFlow . . . . .	52

<b>4</b>	<b>Design of the control system</b>	54
4.1	Outline of the system . . . . .	54
4.1.1	Human-Robot Interaction . . . . .	56
4.1.2	OpenAI Gym Cozmo Environment . . . . .	57
4.1.3	Reinforcement learning algorithms implementations . . . . .	59
<b>5</b>	<b>Experimental results (WIP)</b>	60
<b>6</b>	<b>Conclusions (WIP)</b>	61
<b>A</b>	<b>Reinforcement Learning</b>	62
A.1	Bellman Equation . . . . .	62
A.2	Dynamic programming . . . . .	64

# Chapter 1

## Introduction

TODO (Macaluso P.):

- General discussion about autonomous driving
- General discussion about Reinforcement Learning and recent results with Deep Reinforcement Learning
- The increasing interests in real-world problems and not only simulations
- Focus on the object of the thesis with motivation, description of the procedure followed and results obtained

## Chapter 2

# Reinforcement Learning

Reinforcement learning (RL) is a field of Machine Learning that is experiencing a period of great fervour in the world of research, fomented by recent progress in deep learning (DL) which opened the doors to function approximation developing what is nowadays known as deep reinforcement learning (deep RL). RL represents the third paradigm of machine learning alongside supervised and unsupervised learning. The idea underlying this research field is that the learning process to solve a decision-making problem consists in a sequence of trial and error where the *agent*, the protagonist of RL, could discover and discriminate valuable decisions from penalising ones exploiting information given by a *reward signal*. This interaction has a strong correlation with what human beings and animals do in the real world to forge their behaviour.

Before discussing the results of this thesis, it is good to delineate what today represents the state-of-the-art to understand the universe behind this paradigm better. Indeed, the exploration of this field of research is the main aim of this chapter: the first section begins with the definition of the notation used and with the theoretical foundations behind traditional RL, then in the second section it moves progressively towards what is deep RL through an introduction to the fundamentals of deep learning and a careful discussion of the most essential algorithms paying more attention to those used during the thesis project. The last section aims to illustrate the starting point and ideas of the thesis, drawing what the scenario of deep RL applied to autonomous systems and real-world robotic tasks is today.

The elaboration of this chapter is inspired by [?], [?], [?], [?] and [?].

## 2.1 Fundamentals of reinforcement learning

Reinforcement learning is a computational approach to sequential decision making. It provides a framework that is exploitable with decision-making problems that are unsolvable with a single action and need a sequence of actions, a broader horizon, to

be solved. This section aims to present the fundamental ideas and notions behind this research field in order to help the reader to develop a baseline useful to approach section 2.2 on page 14 about deep reinforcement learning.

### 2.1.1 The reinforcement learning problem

#### Agent, environment and reward

The primary purpose of RL algorithms is to learn how to improve and maximise a future reward by relying on interactions between two main components: the *agent* and the *environment*.

The *agent* is the entity that interacts with the environment by making decisions based on what it can observe from the state of the surrounding situation. The decisions taken by the agent consist of *actions* ( $a_t$ ). The agent has no control over the environment, but actions are the only means by which it can modify and influence the environment. Usually, the agent has a set of actions it can take, which is called *action space*. Some environments have *discrete* action spaces, where only a finite number of moves are available (e.g.  $\mathcal{A} = [\text{North}, \text{South}, \text{East}, \text{West}]$  choosing the direction to take in a bidimensional maze). On the other side, there are *continuous* action spaces where actions are vectors of real values. This distinction is fundamental to choose the right algorithm to use because not all of them could be compatible with both types: according to the needs of the specific case, it may be necessary to modify the algorithm to make it compatible. The sequence of states and actions is named *trajectory* ( $\tau$ ): it is helpful to represent an episode in the RL framework.

The *environment* ( $E$ ) represents all the things that are outside the agent. At every action received by the agent, it emits a *reward* and an *observation* of the environment.

The *reward*  $r_t$  is a scalar feedback signal that defines the objective of the RL problem. This signal allows the agent to be able to distinguish positive actions from negative ones in order to reinforce and improve its behaviour. It is crucial to notice that the reward is local: it describes only the value of the latest action. Furthermore, actions may have long term consequences, delaying the reward. As it happens with human beings' decisions, receiving a conspicuous reward at a specific time step does not exclude the possibility to receive a small reward immediately afterwards and sometimes it may be better to sacrifice immediate reward to gain more rewards later.

In this context, many features make RL different from the other two learning paradigm. Firstly, there is no supervisor: when the agent has to decide what action to take, there is no entity that can tell him what the optimal decision is in that specific moment. The agent receives only a reward signal which may delay compared to the moment in which it has to perform the next action. This fact

brings out another significant difference: the importance of time. The sequentiality links all actions taken by the agent, making resulting data no more independent and identically distributed (i.i.d).

### The concept of return

Given these definitions, it is noticeable that the primary purpose of the agent is to maximise the cumulative reward called *return*.

The *return*  $g_t$  is the total discounted reward starting from timestep  $t$  defined by eq. (2.1) where  $\gamma$  is a *discount factor*.

$$g_t = r_{t+1} + \gamma r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}, \quad \gamma \in [0,1) \quad (2.1)$$

Not only the fact that animal and human behaviour show a preference for immediate rewards rather than for the future ones motivates the presence of this factor, but it is also mathematically necessary: an infinite-horizon sum of rewards may not converge to a finite value. Indeed, the return function is a geometric series, so, if  $\gamma \in [0,1)$ , the series converges to a finite value equal to  $1/(1 - \gamma)$ . For the same convergence sake, the case with  $\gamma = 1$  makes sense only with a finite-horizon cumulative discounted reward.

### States and observations

The other data emitted by the environment is the *observation* ( $o_t$ ) that is related to the *state* ( $s_t$ ). It represents a summary of information that the agent uses to select the next action, while the *state* is a function of the *history* the sequence of observation, actions and rewards at timestep  $t$  as shown in eq. (2.2).

$$h_t = o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t, \quad s_t = f(h_t) \quad (2.2)$$

The state described above is also called *agent state*  $s_t^a$ , while the private state of the environment is called *environment state*  $s_t^e$ . This distinction is useful for distinguishing fully observable environments where  $o_t = s_t^e = s_t^a$ , from partially observable environments where  $s_t^e \neq s_t^a$ . In the first case, the agent can observe the environment state directly, while in the second one, it has access to partial information about the state of the environment. Beyond the fact that this chapter will focus on fully observable environments, the distinction between state and observation is often unclear and, conventionally, the input of the agent is composed by the reward and the state as shown in fig. 2.1 on the next page.

Furthermore, a state is called *informational state* (or *Markov state*) when it contains all data and information about its history. Formally, a state is a Markov state if and only if satisfies eq. (2.3).

$$\mathbb{P}[s_{t+1}|s_t] = \mathbb{P}[s_{t+1}|s_1, \dots, s_t] \quad (2.3)$$



It means that the state contains all data and information the agent needs to know to make decisions: the whole history is not useful anymore because it is inside the state. The environment state  $s_t^e$  is a Markov state.

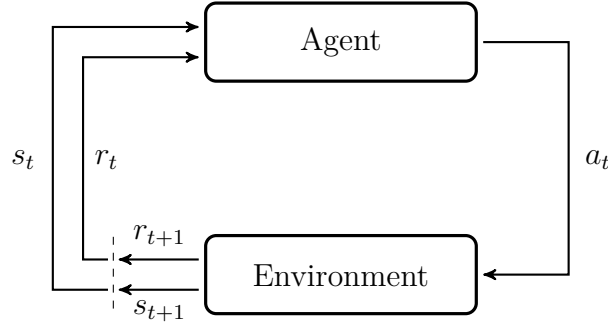


Figure 2.1. Interaction loop between Agent and Environment. The reward and the state resulting from taking an action become the input of the next iteration. [?]

## The Markov decision problem

With all the definitions shown so far, it is possible to formalise the type of problems on which RL can unleash all its features: the Markov decision process (MDP), a mathematic framework to model decision processes. Its main application fields are optimization and dynamic programming.

A MDP is defined by:

$$\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$$

where  $\mathcal{S}$  is a finite set of states

$\mathcal{A}$  a finite set of actions

$\mathcal{P}$  a state transition probability matrix  $\mathcal{P}_{ss'}^a = \mathbb{P}[s_{t+1} = s' | s_t = s, a_t = a]$

$\mathcal{R}$  a reward function  $\mathcal{R}_s^a = \mathbb{E}[r_{t+1} | s_t = s, a_t = a]$

$\gamma$  a discount factor such that  $\gamma \in [0, 1]$

(2.4)

The main goal of an MDP is to select the best action to take, given a state, in order to collect the best reward achievable.

## Policies, models and value functions

In this quick overview of central units of RL, the components that may compose the agent, the brain of the RL problem can not be missing: they are the *model*, the *policy* and the *value function*.

A *model* consist of information about the environment. These data must not be confused with the ones provided by states and observations: they make it possible to infer prior knowledge about the environment, influencing the behaviour of the agent.

A *policy* is the representation of the agent's behaviour. It is a function that describes the mapping from states to actions. The policy is represented by  $\pi$  and it may be deterministic  $a_t = \pi(s_t)$  or stochastic  $\pi(a_t|s_t) = \mathbb{P}[a_t|s_t]$ . In this perspective, it is evident that the central goal of RL is to learn an optimal policy  $\pi^*$ . The optimal policy is a policy which can show to the agent what the most profitable way to achieve the maximum return is, what is the best action to do in a specific situation. In order to learn the nature of the optimal policy, RL exploits value functions.

A *value function* represents what is the expected reward that the agent can presume to collect in the future, starting from the current state. The reward signal represents only a local value of the reward, while the value function provides a broader view of future rewards: it is a sort of prediction of rewards. It is possible to delineate two main value functions: the *state value* function and the *action value* function.

- The *State Value Function*  $V^\pi(s)$  is the expected return starting from the state  $s$  and always acting according to policy  $\pi$ .

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi}[g_t | s_0 = s] \quad (2.5)$$

- The *Action Value Function*  $Q^\pi(s)$  is the expected return starting from the state  $s$ , taking an action  $a$  and then always acting according to policy  $\pi$ .

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[g_t | s_0 = s, a_0 = a] \quad (2.6)$$

### 2.1.2 Bellman equations

Both eqs. (2.5) and (2.6) satisfy recursive relationships between the value of a state and the values of its successor states. It is possible to see this property deriving *Bellman equations* [?] – shown in eq. (2.7) and demonstrated in appendix A.1 on page 62 – where  $s_{t+1} \sim E$  means that the next state is sampled from the environment  $E$  and  $a_{t+1} \sim \pi$  shows that the policy  $\pi$  determines the next action.

$$\begin{aligned} V^\pi(s_t) &= \mathbb{E}_{a_t \sim \pi, s_{t+1} \sim E}[r(s_t, a_t) + \gamma V^\pi(s_{t+1})] \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r|s, a) [r + \gamma V^\pi(s')] \\ Q^\pi(s_t, a_t) &= \mathbb{E}_{s_{t+1} \sim E}[r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi}[Q^\pi(s_{t+1}, a_{t+1})]] \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r|s, a) [r + \gamma Q^\pi(s', a')] \end{aligned} \quad (2.7)$$

$r(s_t, a_t)$  is a placeholder function to represent the reward given the starting state and the action taken. As discussed above, the goal is to find the optimal policy  $\pi^*$  to exploit. It can be done using *Bellman optimality equations* defined in eq. (2.8).

$$\begin{aligned}
 V^*(s_t) &= \max_a \mathbb{E}_{s_{t+1} \sim E} [r(s_t, a) + \gamma V^*(s_{t+1})] \\
 &= \max_a \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma V^*(s')] \\
 Q^*(s_t, a_t) &= \mathbb{E}_{s_{t+1} \sim E} [r(s_t, a_t) + \gamma \max_{a'} Q^*(s_{t+1}, a')] \\
 &= \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma \max_{a'} Q^*(s', a')]
 \end{aligned} \tag{2.8}$$

Therefore, value functions allow defining a partial ordering over policies such that

$$\pi \geq \pi' \text{ if } V_\pi \geq V_{\pi'}, \forall s \in \mathcal{S}$$

This definition is helpful to enounce the *sanity theorem*. It asserts that for any MDP there exists an optimal policy  $\pi^*$  that is better than or equal to all other policies,  $\pi^* \geq \pi, \forall \pi$ , but also that all optimal policies achieve the optimal state value function and the optimal action-value function.

The solution of the Bellman optimality equation is not linear and, in general, there is no closed-form solution. For this reason, there are many iterative methods: sections 2.1.4 and 2.1.5 on page 9 and on page 10 contain some of them.

### 2.1.3 Approaches to Reinforcement Learning

Every agent consists of an RL algorithm that it exploits to maximise the reward it receives from the environment. Every single algorithm has its singularity, and it could work with a specific application field which depends on the particular approach it supports. Understanding differences among these groups is useful to adequately understand what type of algorithm satisfies better the needs of a specific problem. Nowadays, RL algorithms are numerous, and drawing the complete picture behind them could be a complicated purpose. The distinctions presented in this section aims to describe the most crucial distinctions that are useful in the context of the thesis without claiming to be exhaustive.

#### Components of learning

The first worthy distinction between RL algorithms can be made analysing how the algorithms exploit the different components of the agent: indeed it is possible to explain the main strategies in RL using *policy*, *model* and *value function* defined in section 2.1.1 on page 5.

One of the most crucial aspects of an RL algorithm is the question of whether the agent has access to or learns the model of the environment: this element enables

the agent to predict state transitions and rewards. A method is *model-free* when it does not exploit the model of the environment to solve the problem. All the actions made by the agent results from direct observation of the current situation in which the agent is. It takes the observation, does computations on them and then select the best action to take. This last representation is in contrast with *model-based* methods. In this case, the agent tries to build a model of the surrounding environment in order to infer information useful to predict what the next observation or reward would be. Both groups of methods have strong and weak sides. Ordinarily, *model-based* methods show their potential in a deterministic environment (e.g. board game with rules). In these contexts, the presence of the model enables the agent to plan by reasoning ahead, to recognise what would result from a specific decision before taking action. The agent can extract all this knowledge and learn an optimal policy to follow. However, this opportunity is not always achievable: the model may be partially or entirely unavailable, and the agent would have to learn the model from its experience. Learning a model is radically complex and may lead to various hurdles to overcome: for instance, the agent can exploit the bias present in the model, producing an agent which is not able to generalise in real environments. On the other hand, model-free methods tend to be more straightforward to train and tune because it is usually hard to build models of a heterogeneous environment. Furthermore, model-free methods are more popular and have been more extensively developed and tested than model-based methods.

The use of policy or value function as the central part of the method represents another essential distinction between RL algorithms. The approximation of the policy of the agent is the base of *policy-based* methods. The representation of the policy is usually a probability distribution over available actions. This method points to optimise the behaviour of the agent directly and may ask manifold observations from the environment: this fact makes this method not so sample-efficient. On the opposite side, methods could be *value-based*. In this case, the agent is still involved in finding the optimal behaviour to follow, but indirectly. It is not interested anymore about the probability distribution of actions. Its main objective is to determine the value of all actions available, choosing the best one. The main difference from the policy-based method is that this method can benefit from other sources, such as old policy data or replay buffer.

## Learning settings

The learning setting could be *online* or *offline*. In the first case, the learning process is done in parallel or concurrently while the agent continues to gather new information to use, while the second one progresses toward learning using limited data. Generalisation becomes a critical problem in the latter approach because the agent is not able to interact anymore with the environment. In the context of this thesis, what matters is *online learning*: the learning phase is not bound to already

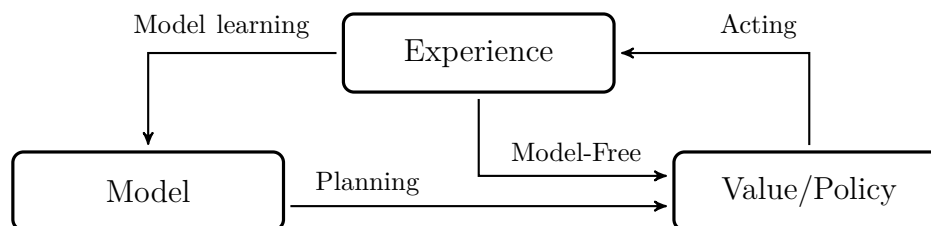


Figure 2.2. Overview of components of an agent with their relation in respect of different approaches of learning. Model-Free methods works with the experience, value functions and the policy, while model-based techniques tries to build up a model to derive value functions and policy to act in the environment.

gathered data, but the whole process goes on using both old data coming from replay buffers and brand new data obtained in the most recent episode.

Another significant difference in RL algorithms consists of the distinctive usage of the policy to learn. *On-policy* algorithms profoundly depend on the training data sampled according to the current policy because they are designed to use only data gathered with the last learned policy. On the other hand, an *off-policy* method can use a different source of valuable data for the learning process instead of direct experience. This feature allows the agent to use, for instance, large experience buffers of past episodes. In this context, these buffers are usually randomly sampled in order to make the data closer to being independent and identically distributed (i.i.d): random extraction guarantees this fact.

### 2.1.4 Dynamic programming

Dynamic programming (DP) is one of the approaches used to resolve RL problems. Formally, it is a general method to explain complex problems by breaking them into more manageable sub-problems. After solving all sub-problems, it is possible to sum them up in order to obtain the final solution to the whole original problem. This technique provides a practical framework to solve MDP problems and to observe what is the best result achievable from it, but it assumes to have full knowledge about the specific problem. For this reason, it applies primarily to model-based problems.

Furthermore, dynamic programming methods *bootstrap*: it means that these strategies use one or more estimated values in the update step for the same kind of estimated value, leading to results more sensitive to initial values.

### Policy Iteration

The *policy iteration* aims to find the optimal policy by directly manipulating the starting policy. However, before proceeding with this process, a proper evaluation

of the current policy is essential. This procedure can be done iteratively following algorithm A.1 on page 63 where  $\theta$  is the parameter that defines the accuracy: the more the value is closer to 0, the more the evaluation would be precise.

*Policy improvement* is the second step towards policy iteration. Intuitively, it is possible to find a more valuable policy than the starting one by changing the action to select in a specific state with a more rewarding one. The key to check if the new policy is better than the previous one is to use the action-value function  $Q_\pi(s, a)$ . This function returns the value of taking action  $a$  in the current state  $s$  and, after that, following the existing policy  $\pi$ . If  $Q_\pi(s, a)$  is higher than  $V_\pi(s)$ , so the action selected is better than the action chosen by the current policy, and consequently, the new policy would be better overall.

Policy improvement theorem is the formalisation of this fact: appendix A.2 on page 64 shows its demonstration. Thanks to this theorem, it is reasonable to act greedily to find a better policy starting from the current one iteratively selecting the action that produces the higher  $Q_\pi(s, a)$  for each state.

The iterative application of policy improvement stops after an improvement step that does not modify the initial policy, returning the optimal policy found.

## Value Iteration

The second approach used by dynamic programming to solve Markov decision processes is *value iteration*. Policy iteration is an iterative technique that alternate evaluation and improvement until it converges to the optimal policy. On the contrary, value iteration uses a modified version of policy evaluation to determine  $V(s)$  and then it calculates the policy. The pseudocode of this method is available algorithm A.2 on page 65.

## Generalised Policy Iteration

Generalised Iteration Policy (GPI) indicates the idea underlying the interaction between evaluation and improvement steps seen in value and policy iteration. Figure 2.3 on the next page reports how the two processes compete and cooperate to find the optimal value function and an optimal policy. The first step, known as policy evaluation step, exploits the current policy to build an approximation of the value function. The second step, known as policy improvement step, tries to improve the policy starting from the current value function. This iterative scheme of dynamic programming can represent almost all reinforcement learning algorithm.

### 2.1.5 Model-free approach

As reported in the previous section, having a comprehensive knowledge of the environment is at the foundation of dynamic programming methods. However,

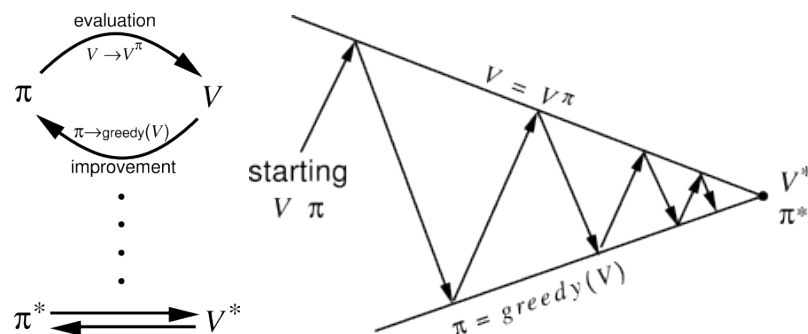


Figure 2.3. Generalised policy iteration schema [?]. Value and policy functions compete and cooperate to reach the joint solution: the optimal value function and an optimal policy.

this fact is not always accurate in practice, where it is infrequent to have a full understanding of how the world works. In these cases, the agent has to infer information using its experience, so it has to exploit model-free methods, based on the assumption that there is no prior knowledge about state transitions and rewards. This section intends to provide a brief description of two model-free approaches to prediction and control: Monte Carlo (MC) methods and Temporal-Difference (TD) ones.

## Monte Carlo learning

Monte Carlo methods [?, Chapter 6] can learn from episodes of experience using the simple idea that averaging sample returns provide the value. This lead to the main caveat of these methods: they work only with episodic MDPs because the episode has to terminate before it is possible to calculate any returns. The total reward accumulated in an episode and the distribution of the visited states is used to calculate the value function while the improvement step is carried out by making the policy greedy concerning the value function.

This approach brings to light the exploration dilemma about how it is possible to guarantee that the algorithm will explore all the states without prior knowledge of the whole environment.  $\epsilon$ -greedy policies are exploited instead of full greedy policy to solve this problem. An  $\epsilon$ -greedy policy is a policy that acts randomly with probability  $\epsilon$  and follows the policy learned with probability  $(1 - \epsilon)$ .

Unfortunately, even though Monte Carlo methods are simple to implement and they are unbiased because they do not bootstrap, they require a high number of iteration to converge. Furthermore, they have a wide variance in their value function estimation due to lots of random decisions within an episode.

## Temporal Difference learning

Temporal Difference (TD) is an approach made combining ideas from both Monte Carlo methods and dynamic programming. TD is a model-free method like MC but uses bootstrapping to make updates as in dynamic programming. The central distinction from MC approaches is that TD methods calculate a *temporal error* instead of using the total accumulated reward. The temporal error is the difference between the new estimate of the value function and the old one. Furthermore, they calculate this error considering the reward received at the current time step and use it to update the value function: this means that these approaches can work with continuing (non-terminating) environments. This type of update reduces the variance compared to Monte Carlo one but increases the bias in the estimate of the value function because of bootstrapping.

The fundamental update equation for the value function is shown in eq. (2.9), where *TD error* and *TD target* are in evidence.

$$V(s_t) \leftarrow V(s_t) + \alpha \underbrace{\left( \overbrace{r_{t+1} + \gamma V(s_{t+1})}^{\text{TD target}} - V(s_t) \right)}_{\text{TD error } (\delta_t)} \quad (2.9)$$

Two TD algorithms for the control problem which are worth quoting because of their extensive use to solve RL problems are *SARSA* (*State-Action-Reward-State-Action*) and *Q-Learning*.

*SARSA* is an on-policy temporal difference algorithm whose first step is to learn an action-value function instead of a state-value function. This approach leads to focus not to estimate the specific value of each state, but to determine the value of transitions and state-action pairs. Equation (2.10) represents the update function of *SARSA*, while algorithm 2.1 on the following page summarise its pseudocode.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)] \quad (2.10)$$

*Q-learning* [?] is an off-policy TD control algorithm which represents one of the early revolution and advance in reinforcement learning. The main difference from SARSA is the update rule for the Q-function: it selects the action in respect of an  $\epsilon$ -greedy policy while the Q-function is refreshed using a greedy policy based on the current Q-function using a max function to select the best action to do in the current state with the current policy. Equation (2.11) represents the update function of *Q-learning*, while algorithm 2.2 on the following page summarise its pseudocode.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (2.11)$$



**Algorithm 2.1:** SARSA (on-policy TD control) for estimating  $Q \approx q_*$ 


---

**Input:** step size  $\alpha \in (0,1]$ , small  $\epsilon > 0$

- 1 Initialise  $Q(s, a) \forall s \in \mathcal{S}, a \in \mathcal{A}$  arbitrarily, except that  $Q(\text{terminal}, \cdot) = 0$
- 2 **foreach** *episode* **do**
- 3     Initialise  $s_t$
- 4     Choose  $a_t$  from  $s_t$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
- 5     **repeat**
- 6         Take action  $a_t \rightarrow$  obtain  $r_{t+1}$  and  $s_{t+1}$
- 7         Choose  $a_{t+1}$  from  $s_{t+1}$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
- 8          $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$
- 9          $s_t \leftarrow s_{t+1} ; a_t \leftarrow a_{t+1}$
- 10     **until**  $s_t$  is terminal
- 11 **end**

---

**Algorithm 2.2:** Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$ 


---

**Input:** step size  $\alpha \in (0,1]$ , small  $\epsilon > 0$

- 1 Initialise  $Q(s, a) \forall s \in \mathcal{S}, a \in \mathcal{A}$  arbitrarily, except that  $Q(\text{terminal}, \cdot) = 0$
- 2 **foreach** *episode* **do**
- 3     Initialise  $s_t$
- 4     Choose  $a_t$  from  $s_t$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
- 5     **repeat**
- 6         Take action  $a_t \rightarrow$  obtain  $r_{t+1}$  and  $s_{t+1}$
- 7         Choose  $a_{t+1}$  from  $s_{t+1}$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
- 8          $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$
- 9          $s_t \leftarrow s_{t+1} ; a_t \leftarrow a_{t+1}$
- 10     **until**  $s_t$  is terminal
- 11 **end**

---

**Temporal Difference Lambda Learning**

As reported previously, Monte Carlo and Temporal Difference learning perform updates in different ways. The first approach exploits the total reward to update the value function, while the second one, on the other hand, works with the reward of the current step. Temporal Difference Lambda, also known as  $\text{TD}(\lambda)$  [?, Chapter 7,12], represents a combination of these two procedures and it takes into account the results of each time step together with the weighted average of those returns. The idea of calculating TD target looking n-steps into the future instead of considering only a single step is the baseline of  $\text{TD}(\lambda)$ . This lead to the formalisation of the

$\lambda$ -weighted return  $G_t^\lambda$  presented in eq. (2.12).

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \quad (2.12)$$

TD( $\lambda$ ) implementation takes into account an additional variable called eligibility trace  $e_t(s_t)$  which indicates how much learning should be carried out for each state for each timestep. It aims to describe how much the agent encountered a specific state recently and eq. (2.13) describes the updating rule of this value where the  $\lambda$  represents the trace-decay parameter.

$$e_t(s) = \gamma \lambda e_{t-1}(s) + \mathbb{1}(s = s_t) \quad (2.13)$$

### 2.1.6 Model-based approach

Heretofore, the focus of this section was on methods which have no prior knowledge of the environment, since this thesis grows on model-free foundations. Despite this point, it is worth to summarise the main concepts behind model-based approaches. Model-based methods gather information to enable the ability of planning, which can enhance the sample efficiency of the algorithm.

There are two primary principles to model-based learning. The first one implies to assemble a model starting from prior knowledge and to exploit it to calculate the policy and the value-function, while the second one is to infer the model from the environment by sampling experience. The central drawback of the first technique is that prior knowledge could be not as accurate as expected, leading to sub-optimal results. Consequently, the preferred way to learn is the second one.

The decisive point behind these approaches is that they are more sample-efficient concerning model-free ones: they require fewer data to learn a policy. On the other hand, the algorithm must learn the policy as well as the model: this translates to two different sources of approximation errors and an increase of computational complexity.

## 2.2 Deep reinforcement learning

The strategies shown so far works smoothly with systems with well-defined states and actions. In this context, it is reasonable to use lookup tables to describe the problem: state-value function  $V$  has an entry for each state while in action-value function  $Q$  has an entry for each state-action pair. It is easy to understand how this setting cannot scale up with very large MDPs: problems regarding the availability of memory arise as it becomes difficult to manage the storage of a large number of states and actions. Also, there may be obstacles concerning the slowness of learning the value of each state individually. Furthermore, the tabular form could lead to

expensive computation in linear lookup and can not work with continuous action and state space.

Function approximators represent the solution to overwhelm this problem. The underlying intention is to use a vector  $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$  to estimate state-value and action-value function as shown in eq. (2.14), generalise from seen states to unseen states and finally update parameter  $\theta$  using MC or TD Learning strategies.

$$\begin{aligned} V(s, \theta) &\approx V_\pi(s) \\ Q(s, a, \theta) &\approx Q_\pi(s, a) \end{aligned} \tag{2.14}$$

In these terms, function approximators can be considered as a mapping from the vector  $\theta$  to the value function. This choice leads to a reduction in the number of parameters to learn and consequently to a system which can generalise better in fewer training samples.

Nowadays, since its widespread use in research, neural networks represent the most intuitive option to take as function approximator: it reduces the training time for high dimensional systems, and it requires less space in memory. This point represents the bridge between traditional reinforcement learning and recent discoveries in the theory of deep learning. Thanks to the last decade great fervour of deep learning, neural networks have become the fundamental tool to exploit as function approximator to develop deep reinforcement learning (Deep RL) which accomplished remarkable results. One of the first steps towards Deep RL and general artificial intelligence – an AI broadly applicable to a different set of various environments – was done by DeepMind with their pioneering paper [?] and the consequent [?].

Because of the nature of this work, the focus of this section will be on model-free algorithms. This section aims to explain the state-of-the-art and the leading theory behind Deep RL framework together with an overview about deep learning and to define two deep actor-critic algorithms used in the experiments of this thesis: Deep Deterministic Policy Gradient (DDPG) and Soft Actor-Critic (SAC).

## 2.2.1 Fundamentals of Deep Learning

### Artificial Neural Networks

Deep learning (DL) is an approach to learning based on a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  parametrised with  $w \in \mathbb{R}^{n_w}$  ( $n_w \in \mathbb{N}$ ) such that  $y = f(x; w)$ .

The starting point of this research field is the artificial neuron, inspired by the biological neuron from the brain of animals and human being. A neuron consists of numerous inputs called *dendrites* coming from preceding neurons. Therefore, the neuron elaborates the input and, only if the value reaches a specific potential, it *fires* through its single output called *axon*.

The neuron elaborates the inputs by taking the weighted sum, adding a bias  $b$  and applying an activation function  $f$  following the relation  $y = f(\sum_n w x_i + b)$ . Figure 2.4 shows the parallel comparison between the biological neuron and the artificial one. The set of parameter  $w$  needs to be adjusted to find a good parameter set: this process is called *learning*.

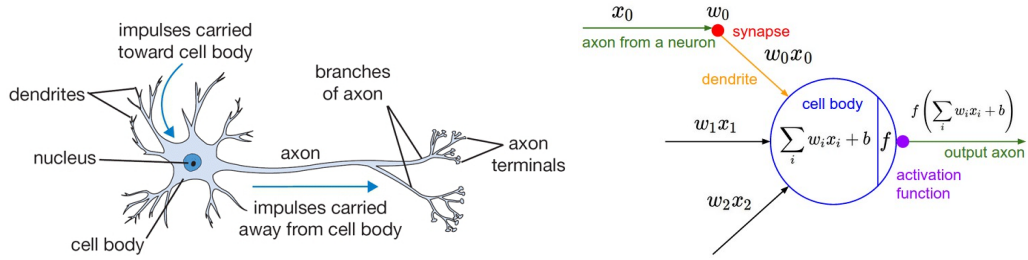


Figure 2.4. Comparison between biological neuron (left) and artificial neuron (right). The artificial neuron designs the dendrites as weighted inputs and returns the sum through an activation function. [?].

A deep neural network (NN) organises a set of artificial neurons in a series of processing layers to which correspond non-linear transformation. The whole sequence of these alterations directs the learning process through different levels of abstraction [?]. To better understand the nature of a deep neural network, it is convenient to describe a neural network with one fully-connected layer represented by fig. 2.5 on the next page.

$$\begin{aligned} h &= g(w_1 \cdot i + b_1) \\ o &= w_2 \cdot h + b_2 \end{aligned} \tag{2.15}$$

The input layer receives as input a column vector of input-features  $i$  of size  $n \in \mathbb{N}$ . Every value of the hidden-layer represented by a vector  $h$  of size  $m \in \mathbb{N}$  is the result of a transformation of the input values given by eq. (2.15) where  $w_1$  is a matrix of size  $m \times n$  and  $b_1$  is a bias term of size  $m$ .  $g$  is a non-linear parametric function called activation function, which represents the core of neural networks. Subsequently, the second and last transformation manipulates the hidden layer  $h$  to produce the values of the output layer following eq. (2.15) using  $w_2$  with size  $o \times m$  and  $b_2$  with size  $o$ .

### Learning process

The learning process aims to seek a set of parameters  $\theta$  that results in the best possible function approximation for a specific objective. In supervised learning, the actual output  $Y$  is available for each particular input  $X$ , and it is used to update the parameters. The learning process can be carried out iteratively according to the following steps.

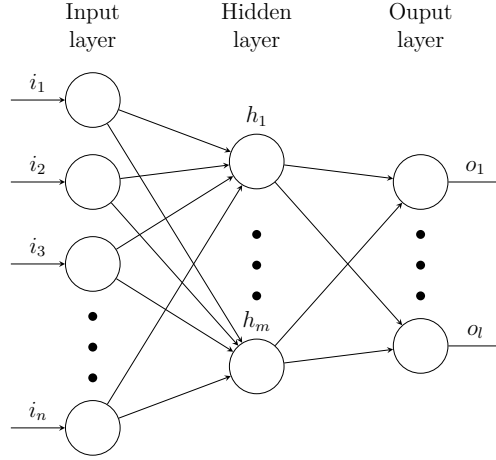


Figure 2.5. An example representation of a deep neural network with one fully-connected hidden layer.

**Forward pass** The input  $X$  is forwarded through the neural network and the output  $Y_{pred} = f(X, \theta)$  is gathered.

**Loss** The resulting predicted value  $Y_{pred}$  is compared with the actual value  $Y$  computing the loss function  $L(\theta)$ . There are a lot of loss functions available to satisfy the particular needs of specific learning tasks. The *error* is the difference between the output of the neural network for a specific input data and the actual value: it is essential to calculate the loss function. One of the most exploited loss function is the *Mean-Squared Error* (MSE) [?] shown in eq. (2.16) which works with L2-distance.

$$L(y, \hat{y}) = (y_{\theta} - y)^2$$

$$J = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) \quad (2.16)$$

**Backpropagation** The next step is the computation of the global gradient of the loss function  $\nabla L(\theta)$ , which is carried out together with its backpropagation through the network. The backpropagation algorithm [?] calculates the local gradient of loss for each neuron in the hidden layers. The concept underlying this procedure and shown in eq. (2.17) is the *chain rule* [?], which computes derivatives of composed functions by multiplying local derivatives.

$$y = g(x), \quad z = f(g(x)), \quad \frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \quad (2.17)$$

Therefore, the chain rule is exploited to propagate the calculated global gradient loss  $\frac{\partial L(\theta)}{\partial \theta}$  back through the network, in the opposite direction of the forward pass. The

procedure calculates the local derivatives during the forward pass, while estimates the local gradient of loss during backpropagation determining the multiplication between the local derivative and the local gradient of the loss of the connected neuron of the next layer: if the neuron has multiple connections neurons, the algorithms adds up all the gradients.

**Update** In this final step consists in the update of the weights of all neurons. There are many ways developed through the years to carry out the update phase, but the most common one is the gradient descent. The objective of the gradient descent is to minimise the loss function by refreshing the internal parameters of the network in the negative direction of the gradient loss: this choice leads the function approximation process closer to the minimum at each iteration. Equation (2.18) describes the update rule presented by the gradient descent where  $\alpha$  is the learning rate. The last-mentioned parameter determines how quickly the algorithm should approach the minimum. A higher learning rate leads to a more significant step towards the minimum, which threatens to overshoot the target.

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J \quad (2.18)$$

Nowadays, the technique applied in the majority of research projects is stochastic gradient descent which combines batch learning [?] and gradient descent, but also its various improved extensions and variants, such as ADAM [?] and AdaGrad [?]: these extensions manage to improve the convergence of SGD thanks to the introduction of adaptive learning rates.

## Regularization

The final aim of the learning process is to obtain a function approximator capable of generalising over data. This fact means that a neural network should show performances on unseen data comparable to the one obtained from training data. For this reason, it is necessary an appropriate trade-off between underfitting and overfitting.

A shallow approximated function and insufficient training data with a lack of diversity are the leading cause to the first situation: the network generalises on the data, but the prediction error is always too high for all data points. The phenomenon of overfitting describes the exact contrary of underfitting. The leading cause is too complex approximation function: this lead to a network which scores an excellent performance on training data, but poorly predicts unseen points.

Regularisation [?, ?] represents an approach to overcome and prevent the problem of overfitting. It works extending the loss function with a regularised term  $\Omega(\theta)$  as shown in eq. (2.19) where  $\lambda$  is the regularisation factor.

$$L'(\theta) = L(\theta, Y, Y_{pred}) + \lambda \Omega(\theta) \quad (2.19)$$

Equations (2.20) and (2.21) show two examples of regularisation terms. The first is  $L^2$ -regularisation which exploits the squared sum of the weights  $\theta$  in order to keep the weights small. The second approach is known as  $L^1$ -regularisation: in this case, large weights are less penalised, but this method leads to a sparser solution.

$$L'(\theta) = L(\theta, Y, Y_{pred}) + \lambda \frac{1}{2} \|\theta\|^2 \quad (2.20)$$

$$L'(\theta) = L(\theta, Y, Y_{pred}) + \lambda \frac{1}{2} \|\theta\| \quad (2.21)$$

### Activation function

Equation (2.22) shows the most common activation functions: in general, *ReLU* achieves better performance over a wide variety of tasks, but usually the selection of the best activation function has to be done starting from all information and requirements of the deep learning model.

$$\begin{aligned} \text{Sigmoid} &\rightarrow g(x) = \frac{1}{1 + e^{-x}} \\ \text{Hyperbolic Tangent} &\rightarrow g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \\ \text{Rectified Linear Unit (ReLU)} &\rightarrow g(x) = \max(0, x) \end{aligned} \quad (2.22)$$

### Batch Learning and Normalization

The basic concept underlying *batch learning* [?] is to process a set of  $n$  training samples called also *mini-batches* in the place of a single one. This method works with the gradient averaged over all the samples in the mini-batch: it leads to a more accurate gradient reducing its variance and the training time.

*Batch normalisation* consists of zero-centring and rescaling all data in a specific batch, resulting in a mean of normalised data close to 0 and a variance close to 1. The algorithm presented in algorithm 2.3 on the following page is the one provided in [?]. This method computes the mean  $\mu_\beta$  and the variance  $\sigma_\beta$  element-wise for each spatial position in the batch:  $\epsilon > 0$  is a small value to avoid the division by zero. The batch normalisation layer then processes the resulting normalised value:  $\gamma$  and  $\beta$  are the parameters of this layer that are in addition to the original parameter set  $\theta$  of the Neural Network in the learning process. The new learning dynamic provided by the addition of this class of layer increases the network expressivity: applying this method to the input data and the output of any hidden layer results in the reduction of the training time, better regularisation during learning and the reduction of the overfitting phenomenon.

---

**Algorithm 2.3:** Batch normalisation
 

---

**Input:** Mini-batch  $\mathcal{B} = x_{1\dots m}$ ;  $\gamma$  and  $\beta$  parameters to be learned

**Output:**  $y_i = \text{BN}_{\gamma,\beta}(x_i)$

- 1  $\mu_\beta \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  // Mini-batch mean
  - 2  $\sigma_\beta^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_\beta)^2$  // Mini-batch variance
  - 3  $\hat{x}_i \leftarrow \frac{x_i - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}}$  // Normalisation
  - 4  $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i)$  // Scale and shift
- 

## Convolutional Neural Networks

Sensory reception represents how humans and animals react to changes: it consists of sensors which process the input data and are sensitive to specific stimuli. This system inspires the architecture underlying Convolutional Neural Networks: they could handle efficiently significant input data with many applications in computer vision. Figure 2.6 displays the *LeNet-5* [?] which is capable to recognise digits in images. It represents a perfect example of a standard convolutional neural network architecture: it consists of a series of convolutional layers followed by a subsampling pooling layer. At the end of the convolutional stack, the values map into final hidden layers of the network to compute the final low-dimensional output of the network: fully-connected layers usually compose these final layers. It is possible to suppose that the first layers have to learn low-level features of the input data while succeeding layers are responsible for combining the last-mentioned features in high-level ones.

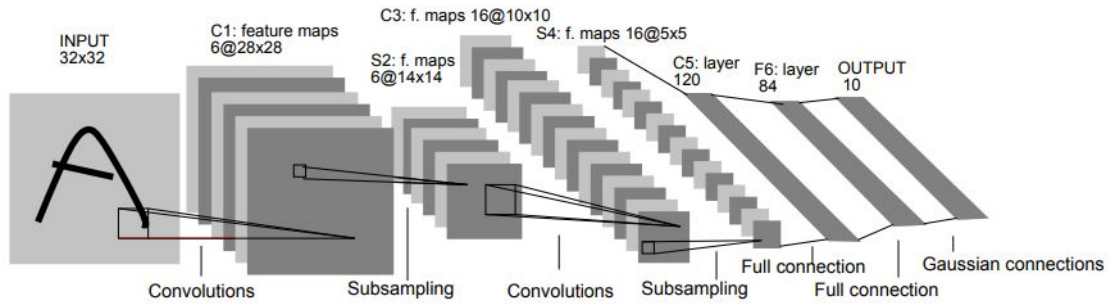


Figure 2.6. Comparison between biological neuron (left) and artificial neuron (right). The artificial neuron designs the dendrites as weighted inputs and returns the sum through an activation function. [?].



**Convolutional Layer** Convolutional layers [?] operates with a set of learnable filters (called kernels) with dimensions  $n \times m$  smaller than the whole input image. The convolution is the basic operation employed in these class of layer: it consists in convolving each filter across the width and height of the input data and computing dot products between the values in the filter and the ones in the input at any position. The result of this operation is a 2-dimensional activation map that contains the results of that filter at every spatial position. In this context, the network can learn filters that detect specific features in the image such as edges, textures and patterns [?].

It is possible to compute the output size  $W_2 \times H_2 \times D_2$  of a pooling layer using starting from the input size  $W_1 \times H_1 \times D_1$  and from the hyperparameters of this class of layer: the number of filters  $K$ , their spatial extent  $F$ , the stride  $S$  and the amount of zero padding  $P$ . The resulting volume size can be calculated using the relations reported in eq. (2.23).

$$\begin{aligned} W_2 &= (W_1 - F + 2P)/S + 1 \\ H_2 &= (H_1 - F + 2P)/S + 1 \\ D_2 &= K \end{aligned} \tag{2.23}$$

The number of parameters introduced by a single kernel is equal to  $F \cdot F \cdot D_1$ , so the convolutional layer has a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  bias.

In a convolutional layer, the number of weights is kept small and then the computation is more efficient than the one of a fully-connected layer: small filters need fewer parameters and less work in the convolutional operation. Besides this motivation, the filters are kept small also because it makes them capable of learning small and low-level features. The great innovation behind convolutional layers is that the same neuron can recognise the related learned features even if they emerge in different locations of the image: this is the property of translation invariance of feature detection with convolutional layers.

**Pooling Layer** It is common to insert a pooling layer in-between successive convolutional layers. The main objective of this class of layers is to apply a down-sampling filter on the input: it progressively reduces the spatial size of the representation, decreasing the number of parameters and the computational cost in the network. It is also useful to control overfitting.

It is possible to compute the size of the output  $W_2 \times H_2 \times D_2$  of a pooling layer using starting from the size of the input  $W_1 \times H_1 \times D_1$  and from the hyperparameters of this class of layer: the spatial extent of the filter  $F$ , the stride  $S$ . The resulting volume size can be calculated using the relations reported in eq. (2.24).

$$\begin{aligned} W_2 &= (W_1 - F)/S + 1 \\ H_2 &= (H_1 - F)/S + 1 \\ D_2 &= D_1 \end{aligned} \tag{2.24}$$

The most common types of pooling layer are the *max-pooling* and the *average-pooling* layer. Both classes return a single value for each position of the filter: the first returns the maximum value, while the second returns the average among the values in the specific section of the input.

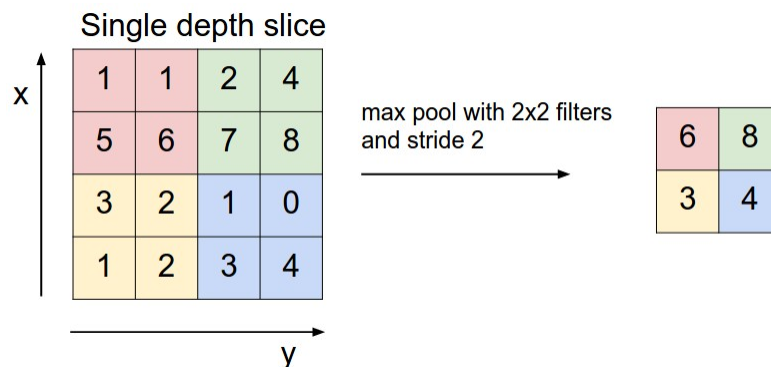


Figure 2.7. Example of a max-pooling operation with a stride of 2 and a spatial extent of 2. [?].

It is worthy of using a pooling layer in a situation where the exact feature position is not relevant but rather whether a particular feature exists in the input at all.

## 2.2.2 Value-based methods

The first class of algorithms to explore is value-based one. They work learning an approximator  $Q_\theta(s, a)$  to infer the optimal action-value function  $Q^*(s, a)$  using an objective function based on Bellman equations. The preponderance of optimisations belonging to this category is *off-policy*: this means that the optimisation step is done using all data collected during the whole training, not only with the most recent policy available. In this configuration, the information used for the learning phase could also come from exploration decision, apart from ones obtained with the most recent policy. Indeed, value-based approaches are more sample efficient because they can reuse data more efficiently, but they are considered less stable than policy gradient ones. Thanks to the relation expressed by  $a(s) = \operatorname{argmax}_a Q_\theta(s, a)$ , it is possible to obtain the policy learned so far from the current action-value function.

### Deep Q-Network (DQN)

This algorithm grows from the ideas underlying Q-Learning [?] shown previously in section 2.1.5 on page 12. The team of DeepMind introduced the Deep Q-Network (DQN) in [?] and in the next cutting-edge paper [?]: they managed to create an algorithm capable of learning to play ATARI video-games online using raw image

and pixels. It works with neural networks as a function approximator, employing convolutional layers in the first layers of the neural network and performing the optimisation with a variant of stochastic gradient descent called RMSprop [?]. The exploited neural network provides as output a probability distribution over all possible discrete actions to determine what is the best action to take.

To overcome the instability problem of value-based methods, DQN utilises two heuristics to narrow instabilities.

**Target Network** The presence of a second network, also called *target network*, enriches the update phase of this algorithm. Equation (2.25) defines the loss function of DQN where  $y_i$  value is computed using the target network instead of the local network. Therefore, the parameters of the target network are hard updated every  $I \in \mathbb{N}$  iterations: this choice precludes instabilities and avoids divergence because of target networks parameters remain fixed for  $I$  iterations.

$$\begin{aligned} L_i(\theta_i) &= \mathbb{E}_{s,a \sim \pi} [(y_i - Q(s, a; \theta_i))^2] \\ y_i &= \mathbb{E}_{s' \sim E} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a] \end{aligned} \quad (2.25)$$

**Experience Memory Replay** Another crucial introduction in this algorithm is *experience memory replay buffer* [?]. Trajectories sampled from the environment are temporally correlated, and this could lead to overfitting the parameters of the neural network because the data are not independent and identically distributed (*i.i.d.*). The setting of this algorithm is *online* because the replay buffer stores  $N_{\text{replay}} \in \mathbb{N}$  replacing old steps as new ones arrive. The experience is collected as tuples  $(s_t, a_t, r_t, s_{t+1})$  using the  $\epsilon$ -greedy policy. The learning phase samples a set of limited tuples called *mini-batch* allowing a wider set of state-action pair in the update of the network and improving the procedure in terms of variance in respect of single tuple update. Trajectories sampled from the environment are temporally correlated, and this could lead to overfitting the parameters of the neural network. Using a batch sample from the replay buffer makes the data *i.i.d.* and consequently improves the learning.

## Improvements of DQN

Further investigation and speculation followed the publication of the thriving DQN. Summing up and comparing these new approaches to original DQN is the main aim of *Rainbow* [?]: it also introduces an algorithm called *Rainbow DQN* with all the techniques proposed. The following paragraphs will delineate three main improvements of DQN.

**Double DQN** The double DQN [?, ?] improvement can handle the intricacy of overestimation of Q-values caused by the maximisation step in eq. (2.25) on the

preceding page. It works with two separate Q-Network with parameters  $\theta$  and  $\theta^-$  for estimating TD-target. It allows for removing the positive bias in estimating action values, leading to less overestimation of Q-learning values, improved stability and performance. In this context, the target  $y_i$  is replaced by eq. (2.26).

$$y_i = \mathbb{E}_{s' \sim E}[r + \gamma Q(s', \underset{a}{\operatorname{argmax}} Q(s', a; \theta_{i-1}); \theta_{i-1}^-) | s, a] \quad (2.26)$$

**Prioritised Experience Replay** The fundamental idea underlying *prioritised experience replay* [?] is precisely to prioritise experiences that contain more crucial information than other ones. An additional value that defines the priority of a specific transition joins each tuple stored in the replay buffer: thanks to this approach, experiences with higher priority has a higher sampling probability and are more likely to remain longer in the replay buffer. It is possible to use *TD-error* to measure the importance of each tuple in the experience. A high TD-error means that the agent behaved better or worse than expected in that particular moment, and therefore, it can learn more from that specific passage.

**Dueling DQN** The dueling DQN architecture [?] presented in fig. 2.8 works decoupling the Q-value estimation in two distinct sequences of fully-connected layers right after convolutional layers.

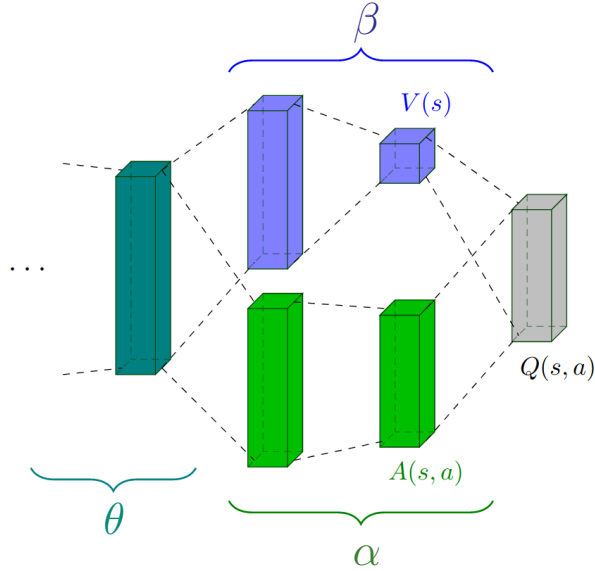


Figure 2.8. Dueling DQN architecture [?]: It consists in two stream to estimate state-value parametrised with  $\beta$  and advantages values parametrised with  $\alpha$  for each action. The last layer represent the combination of these two types of values to obtain the Q-function. [?]

These streams are capable of providing separate estimates of the state-value  $V(s; \theta, \beta)$  and advantage  $A(s, a; \theta, \alpha)$  functions exploited in the end to obtain the Q-value function  $Q(s, a; \theta, \alpha, \beta)$  estimate where  $\theta$  represents the parameters of convolutional layers while  $\alpha$  and  $\beta$  the ones of state-value and advantage function respectively.

The first formalisation of the Q-value function is shown by eq. (2.27), but [?] also suggests a different approach shown in eq. (2.28) with increased stability in practice.

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + \left( A(s, a; \theta, \alpha) - \max_{a' \in \mathcal{A}} A(s, a'; \theta, \alpha) \right) \quad (2.27)$$

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + \left( A(s, a; \theta, \alpha) - \frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} A(s, a'; \theta, \alpha) \right) \quad (2.28)$$

### 2.2.3 Policy gradient methods

Policy gradient algorithms aim to optimise the policy's performance measure in eq. (2.29) by finding a suitable policy  $\pi_\theta(s|a)$  capable of generating a trajectory  $\tau$  that maximises the expected rewards eq. (2.30) instead of learning a value function. Indeed the objective function in policy gradient methods consists of maximising  $J(\theta)$  value by finding a proper policy updating  $\theta$  parameters directly.

$$J(\theta) = \mathbb{E} \left[ \sum_{t=0}^N r(s_t, a_t); \pi_\theta \right] = \sum_{\tau} P(\tau; \theta) r(\tau) \quad (2.29)$$

$$\theta^* = \operatorname{argmax}_{\theta} J(\theta) \quad (2.30)$$

Stochastic gradient ascent is used to refresh the parameters of the policy  $\theta$ . Gradient ascent is the inverse of gradient descent and updates the parameters  $\theta_t$  in the positive direction of the gradient of the policy's performance measure  $\nabla_{\theta} J(\theta)$  following eq. (2.31) where  $\alpha$  is the learning rate which defines the strength of the steps in the direction of the gradient.

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_{\theta} J(\theta_t) \quad (2.31)$$

The main advantage of policy gradient approaches consists in the stability of their convergence: these methods work updating their policy directly at each time step instead of renewing value function from which to derive the policy like value-based methods. Last-mentioned approaches can lead to a radical change in the policy output even for a small change in the value function: this event can cause prominent oscillation during training. Furthermore, policy gradient algorithms can face infinite and continuous action space because the agent estimates the action directly instead of calculating the Q-value for each possible discrete action. The third feature is their ability to learn stochastic policies, useful in uncertain contexts or partially

observable environments. Despite the presence of the advantages just mentioned, policy gradient methods have a substantial disadvantage: they tend to converge to a local maximum instead of the global optimum.

### Actor-Critic Architecture

Actor-critic architecture, shown in fig. 2.9, represents the point of contact between value-based approaches and policy gradient methods. They are policy gradient methods basically but exploit value-function to learn the parameters  $\theta$  of the policy. As its name suggests, these approaches work with two different parts called *actor* and *critic* [?]. The *actor* relates to the policy, while the *critic* deals with the estimation of a value function (e.g. Q-value function). In the context of deep reinforcement learning, they can be represented using neural networks function approximator [?]: the actor exploits gradients derived from the policy gradient theorem and adjusts the policy parameters, while the critic estimates the approximate value function for the current policy  $\pi$ .

Standard practice is to update both networks with the *TD-Error*, discussed in section 2.1.5 on page 12. Estimation made by the critic is useful to determine the contribution that expected values of the current and next state gives to the TD-error. Essentially, the output of the critic contributes to the update of the actor.

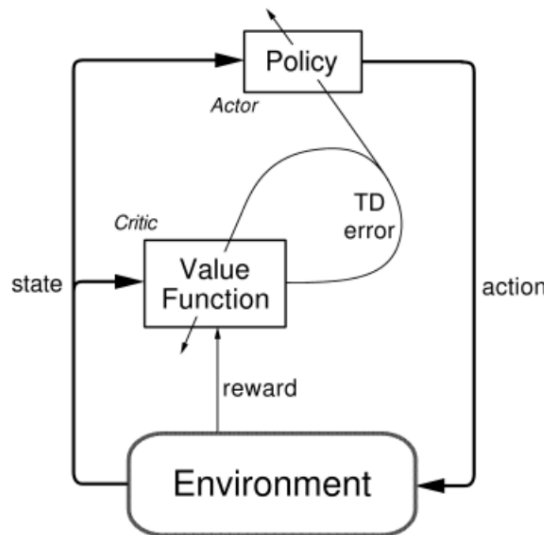


Figure 2.9. Actor-critic architecture schema: the actor represents the policy and maps the input state to an output action, while the critic represents the value function. Both networks can be updated with, e.g. the TD-error, using the contribution of the critic. It is noticeable that the actor uses the critic during the learning process [?].

## 2.2.4 Deep Deterministic Policy Gradient (DDPG)

Deep Deterministic Policy Gradient (DDPG) [?] is a policy gradient algorithm that works learning a Q-function and a policy, and that grows from the deterministic policy gradient algorithm (DPG) [?]. It is a model-free, off-policy, actor-critic algorithm which utilises deep function approximators to learn policies in high-dimensional, continuous action spaces. It can be applied to situations that can not be solved using *DQN* algorithm [?] because of the presence of continuous action spaces. A fine discretisation of the action space to adapt the situation to DQN would lead to an explosion in the number of discrete actions and the curse of dimensionality. *Bellman equation* and *Q-learning* are integral parts of this algorithm. The algorithm concurrently learns a Q-value function and a policy: it uses off-policy data and the Bellman equation to learn the Q-value function and uses the Q-value function to learn the policy.

Usually, in Reinforcement Learning, if the optimal action-value function  $Q^*(s, a)$  is known, then in any given state, the optimal action  $a^*(s)$  can be found by solving eq. (2.32).

$$a^*(s) = \arg \max_a Q^*(s, a) \quad (2.32)$$

When the number of discrete actions is finite, calculating the max poses no problem, because the Q-values can be calculated for each action separately, then directly compared. However, when the action space is continuous, this process becomes highly non-trivial: it would need to be run at every step of the episode, whenever the agent wants to take any action in the environment, and this can not work.

Because the action space is continuous, the function  $Q^*(s, a)$  is presumed to be differentiable concerning the action argument. For this reason, an efficient, gradient-based learning rule for a policy  $\pi(s)$  which exploits that fact can be set up, approximating it with  $\max_a Q(s, a) \approx Q(s, \pi(s))$ .

### Target Networks

DDPG algorithm exploits 4 neural networks: the *local actor*, the *local critic*, the *target actor* and the *target critic*. Actor networks aim is to approximate the policy using parameters  $\theta$  while critic networks approximate the Q-Value function using parameters  $\phi$ .

Initially, actor and critic networks have both randomly initialised parameters. Then the local actor – the current policy – starts to propose actions to the agent, given the current state, starting to populate the experience replay buffer.

When the replay buffer is big enough, the algorithm starts to sample randomly a mini-batch of experiences for each timestep  $t$ . This mini-batch is used to update the local critic minimising the Mean Squared Error (MSE) between the local Q-value and the target one shown in eq. (2.33) on the next page where  $y_i = y_t$  given by eq. (2.36) on the following page and to update the actor policy using the sampled

policy gradient defined in eq. (2.39) on the next page.

$$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \phi))^2 \quad (2.33)$$

We can imagine the target networks as the *labels* of supervised learning.

Also the target networks are updated in this *learning step*. A mere copy of the local weights is not an efficient solution, because it is prone to divergence. For this reason, a *soft* target updates is used. It is given by eq. (2.34) where  $\tau \ll 1$ .

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta \quad (2.34)$$

## Learning Equations

The two fundamental functions in Reinforcement Learning exploited in DDPG are the *action-value function* – see eq. (2.6) on page 6 – and the correspondent *Bellman equation* – see eq. (2.7) on page 6. If this policy is deterministic we can describe it as a function  $\pi : \mathcal{S} \leftarrow \mathcal{A}$  obtaining eq. (2.35) which depends only on the environment.

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E}[r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, \pi(s_{t+1}))] \quad (2.35)$$

This means that it is possible to learn  $Q^\pi$  off-policy, using transition generated by a different stochastic behaviour policy  $\beta$ .

Focusing more and more on DDPG, the Bellman equation is the starting point for learning an approximator to  $Q^*(s, a)$  of Q-Learning. The approximator is parametrised by  $\phi$  and the value network is updated and optimised by minimising the loss defined in eq. (2.36) where  $d_t$  is a flag which indicates whether the state  $s_{t+1}$  is terminal.

$$\begin{aligned} L(\phi) &= \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta, r_t \sim E}[(Q(s_t, a_t | \phi) - y_t)^2] \\ y_t &= r(s_t, a_t) + \gamma(1 - d_t)Q(s_{t+1}, \pi(s_{t+1}) | \phi) \end{aligned} \quad (2.36)$$

It is clear from the eq. (2.36) that the loss is calculated starting from the transitions generated by the policy  $\beta$ . For this reason a great importance in this algorithm is given to *replay buffer* – see section 2.2.2 on page 23 – and *target networks* – see section 2.2.4 on the previous page.

From the policy perspective, the objective is to maximise eq. (2.37) calculating the policy loss through the derivative of the objective function with respect to the policy parameter eq. (2.38). However, since the algorithm is updating the policy in an off-policy way with batches of experience, it is possible to use the mean of the sum of gradients calculated from the mini-batch eq. (2.39) on the following page.

$$J(\theta) = \mathbb{E}[Q(s, a) |_{s=s_t, a=\pi(s_t)}] \quad (2.37)$$

$$\nabla_\theta J(\theta) \approx \nabla_a Q(s, a) \nabla_\theta \pi(s | \theta) \quad (2.38)$$



$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \left[ \nabla_a Q(s, a | \phi) |_{s=s_i, a=\pi(s_i)} \nabla_{\theta} \pi(s | \theta) |_{s=s_i} \right] \quad (2.39)$$

Algorithm 2.4 on the next page shows the pseudocode of the DDPG algorithm.

### Exploration vs. Exploitation

In reinforcement learning for discrete action spaces, exploration is done selecting a random action (e.g. epsilon-greedy). For continuous action spaces, exploration is done adding noise to the action itself. In [?], the authors use Ornstein-Uhlenbeck process [?] to add noise to the action output  $a_t = \pi(s_t | \theta) + \mathcal{N}$ . After that the action is clipped in the correct range.

### Hyperparameters

**Exploration noise** The exploration noise consists of two sets of parameters. The first set refers to the Ornstein-Uhlenbeck Process noise parameters  $\pi, \sigma, \theta$  as reported in [?]. The second one consists of the parameters of  $\epsilon$ , a small value to decrease the impact of the noise on the action. Equation (2.40) describes how the impact of the noise decreases in function of the current episode number  $e$ , where  $\epsilon_{\text{start}}$  represent the starting value of the noise and  $\epsilon_{\text{end}}$  the final one.

$$\epsilon = \epsilon_{\text{start}} - (\epsilon_{\text{start}} - \epsilon_{\text{end}}) \min \left( 1.0, \frac{e}{\epsilon_{\text{decay}}} \right) \quad (2.40)$$

**Replay buffer** The parameters available for replay memory are its maximum size, its minimum size to start the learning phase and the mini-batch size to sample for each learning step.

**Neural Network** The neural network can be considered as a whole complex hyperparameter because it is possible to select among different layers to exploit for specific problems – e.g. the number of layers, the type of layers, the number of hidden features. Given the network architecture, the primary neural networks hyperparameters are the learning rate  $\alpha$  and the update method – e.g. Adaptive Momentum Estimation (ADAM).

**Learning Update** The parameters of the learning phase of the algorithm are mainly two. The first one is  $\gamma$ , the main parameter in the reinforcement learning framework, which characterises the discounted return. The second is the soft target update parameter  $\tau$ , which determinates the entity of the update of the network at each learning step.

**Algorithm 2.4:** DDPG Algorithm [?]

---

**Input:** Initial critic network parameter  $\bar{\phi}$  and actor network parameter  $\bar{\theta}$

- 1 Initialise target network weights  $\bar{\phi} \leftarrow \phi, \bar{\theta} \leftarrow \theta$
- 2 Initialise a random process  $\mathcal{N}$  for action exploration
- 3 **for**  $episode = 1, M$  **do**
- 4     Receive the initial observation state  $s_t \leftarrow s_1$
- 5     **repeat**
- 6         Select action  $a_t = \pi(s_t|\theta) + \mathcal{N}$  and *clip* results
- 7         Execute action  $a_t$  and obtain tuple  $(r_t, s_{t+1}, d_t)$
- 8         Store transition  $(s_t, a_t, r_t, s_{t+1}, d_t)$  in  $\mathcal{D}$
- 9         **if** *it is time to update* **then**
- 10             Sample random minibatch of  $N$  transitions  $(s_t, a_t, r_t, s_{t+1}, d_t)$  from  $\mathcal{D}$
- 11             Compute the target  $y_i = r_i + \gamma(1 - d_i)Q'(s_{i+1}, \pi'(s_{i+1}|\theta)|\phi)$
- 12              $\phi \leftarrow \phi - \lambda_Q \nabla_{\phi} J_Q(\phi)$  // Update critic
- 13              $\theta \leftarrow \theta - \lambda_{\pi} \nabla_{\theta} J_{\pi}(\theta)$  // Update actor
- 14              $\phi \leftarrow \tau\phi + (1 - \tau)\bar{\phi}$  // Soft update critic
- 15              $\theta \leftarrow \tau\theta + (1 - \tau)\bar{\theta}$  // Soft update actor
- 16         **end**
- 17          $s_t \leftarrow s_{t+1}$
- 18     **until**  $s_t$  is terminal
- 19 **end**

**Output:** Optimised parameters  $\theta$  and  $\phi$

---

### 2.2.5 Soft Actor-Critic (SAC)

Soft Actor-Critic (SAC) [?, ?] combines the off-policy actor-critic setup with a stochastic policy (actor), devising a bridge between stochastic policy optimization and DDPG-style approaches. As DDPG, SAC can work in situations characterised by the presence of continuous action spaces, and it is a model-free, off-policy and actor-critic algorithm.

SAC algorithm can overcome some of the problems of DDPG. The latter can achieve excellent performance, but the interaction between the deterministic actor-network and the Q-function makes it difficult to stabilise and brittle concerning hyperparameters and other kinds of tuning [?, ?]. The learned Q-function begins to dramatically overestimate Q-values, which then leads to the policy breaking because it exploits the errors in the Q-function. For this reason, SAC exploits *Clipped Double-Q Learning* also used by Twin Delayed DDPG (TD3) [?]. It learns two Q-functions instead of one and uses the smaller of the two Q-values to form the targets in the Bellman error loss functions.

Another feature of SAC is *entropy regularization* [?, ?, ?, ?]. The policy is trained to maximise a trade-off between expected return and entropy, a measure of

randomness in the policy. This peculiarity is strongly related to the exploration-exploitation trade-off: increasing entropy results in more exploration, which can accelerate learning later on, but it can also prevent the policy from prematurely converging to a local optimum.

### Target Networks

SAC algorithm exploits 5 neural networks: the local stochastic policy network with parameter  $\theta$ , two local Q-Networks with parameters  $\phi_1, \phi_2$  respectively, two target Q-Networks with parameters  $\bar{\phi}_1$  and  $\bar{\phi}_2$  respectively. Their behaviour is the same as the one of DDPG target network: the algorithm updates the target networks following eq. (2.34) on page 28.

Algorithm 2.5 on page 34 shows the pseudocode of the SAC algorithm.

### Entropy-Regularised Reinforcement Learning

*Entropy* represents the average rate at which a stochastic source of data produces information. It is, in simple terms, a quantity which describes how random a random variable is. The motivation behind the use of entropy is that when the data source produces a low-probability value, the event carries more information than when the source data produces a high-probability value.

Let  $x$  be a random variable with probability mass or density function  $P$ . The entropy  $\mathcal{H}$  of  $x$  is computed from its distribution  $P$  according to eq. (2.41).

$$\mathcal{H}(P) = \mathbb{E}_{x \sim P}[-\log P(x)] \quad (2.41)$$

In *entropy-regularised* reinforcement learning the standard objective is generalised by augmenting it with entropy. The agent gets a bonus reward at each time step proportional to the entropy of the policy at that timestep. Assuming an infinite-horizon discounted setting, this changes the RL problem as shown in eq. (2.42) where  $\alpha > 0$  is the temperature parameter that determines the relative importance of the entropy term controlling the stochasticity of the optimal policy.

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \left( R(s_t, a_t, s_{t+1}) + \alpha \mathcal{H}(\pi(\cdot | s_t)) \right) \right] \quad (2.42)$$

It is clear that the standard maximum expected return can be retrieved in the limit as  $\alpha \rightarrow 0$ .

From eq. (2.42) it is possible to derive *state-value function*  $V^{\pi}(s)$  and *action-value function*  $Q^{\pi}(s, a)$  as shown in eq. (2.43) and eq. (2.44) on the next page.

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \left( R(s_t, a_t, s_{t+1}) + \alpha \mathcal{H}(\pi(\cdot | s_t)) \right) \middle| s_0 = s \right] \quad (2.43)$$

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1}) + \alpha \sum_{t=1}^{\infty} \gamma^t \mathcal{H}(\pi(\cdot|s_t)) \middle| s_0 = s, a_0 = a \right] \quad (2.44)$$

From these equations is possible to derive the connection between state-value and action-value function given by eq. (2.45) and the Bellman equation given by eq. (2.46).

$$V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] + \alpha \mathcal{H}(\pi(\cdot|s)) \quad (2.45)$$

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}_{s' \sim P, a' \sim \pi} [R(s, a, s') + \gamma(Q^\pi(s', a') + \alpha \mathcal{H}(\pi(\cdot|s')))] \\ &= \mathbb{E}_{s' \sim P} [R(s, a, s') + \gamma V^\pi(s')] \end{aligned} \quad (2.46)$$

## Learning Equations

SAC algorithm learns a policy  $\pi_\theta$  with  $\theta$  parameter set and two Q-functions  $Q_{\phi_1}, Q_{\phi_2}$  with  $\phi_1$  and  $\phi_2$  parameter sets respectively. The state-value function is implicitly parametrised through the soft Q-function parameters thanks to eq. (2.47). In [?] a function approximator for this function was introduced, but later [?] the authors found it to be unnecessary.

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] + \alpha \mathcal{H}(\pi(\cdot|s)) \\ &= \mathbb{E}_{a \sim \pi} [Q^\pi(s, a) - \alpha \log \pi(a|s)] \\ &\approx Q^\pi(s, \tilde{a}) - \alpha \log \pi(\tilde{a}|s), \quad \tilde{a} \sim \pi(\cdot|s). \end{aligned} \quad (2.47)$$

**Learning Q** Q-functions are learned by Mean Squared Bellman Error (MSBE) minimization, using a target value network to form the Bellman backups using eq. (2.48). Equation (2.47) implicitly parametrise the state-value function.

$$J_Q(\phi_i) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_{\phi_i}(s_t, a_t) - \left( r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} V_{\phi_i^-}(s_{t+1}) \right) \right)^2 \right] \quad (2.48)$$

The update shown makes use of target soft Q-function with parameters  $\phi_i$ , which are calculated, like in the DDPG algorithm, as an exponentially moving average of the soft Q-function parameters [?]. It can be optimised using stochastic gradients.

**Learning the Policy** It is possible to derive SAC starting from the definition of soft policy iteration demonstrated in [?, Section 4]. In particular, the policy has to be learned starting from the minimization of the expected KL-divergence [?, ?] and exploiting eq. (2.49).

$$J_\pi(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} [\mathbb{E}_{a_t \sim \pi_\theta} [\alpha \log \pi_\theta(a_t|s_t) - Q_\phi(s_t, a_t)]] \quad (2.49)$$

As [?] reports, there are several options for the minimization of  $J_\pi(\theta)$ , but the most straightforward one using neural network as function approximator is to apply the

*reparametrization trick*. It works reparametrising the policy using a neural network transformation following eq. (2.50) where  $\epsilon_t$  is a noise vector sampled from some fixed distribution.

$$a_t = f_\theta(\epsilon_t; s_t) \quad (2.50)$$

Therefore, it is possible to rewrite the expectation over actions in eq. (2.49) on the previous page into an expectation over noise as in eq. (2.51).

$$J_\pi(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}}[\alpha \log \pi_\theta(f_\theta(\epsilon_t; s_t) | s_t) - Q_\phi(s_t, f_\theta(\epsilon_t; s_t))] \quad (2.51)$$

To get the policy loss, the final step is to substitute  $Q_\phi$  with one of our function approximators: the choice falls  $\min_{i=1,2} Q_{\phi_i}(s_t, f_\theta(\epsilon_t; s_t))$ , as suggested from the authors of [?].

## Exploration vs. Exploitation

SAC algorithm trains a stochastic policy using *entropy regularization*.  $\alpha$  is the entropy regularisation coefficient which is the parameter that explicitly controls the exploration-exploitation trade-off. A higher  $\alpha$  corresponds to more exploration, while a lower  $\alpha$  corresponds to more exploitation.

This parameter has fundamental importance in the algorithm, and it may vary from environment to environment. Choosing the optimal  $\alpha$  parameter is a non-trivial task that could require careful tuning in order to find the one which leads to the stablest and highest-reward learning. In [?, Section 5], the authors formulated a different maximum entropy reinforcement learning algorithm to overcome this problem. Forcing the entropy to a fixed value is a weak solution because the policy should be free to explore more where the optimal action is uncertain, and to exploit the learned mapping in states with a more clear optimal action. The gradients are computed using eq. (2.52)

$$J(\alpha) = \mathbb{E}_{a_t \sim \pi}[-\alpha \log \pi(a_t | s_t) - \alpha \bar{\mathcal{H}}] \quad (2.52)$$

During the test phase, the algorithm uses the mean action instead of a sample from the distribution learned. This choice tends to improve performance over the original stochastic policy, allowing to see how well the policy exploits what it has learned.

## Hyperparameters

**Entropy regularisation parameter** The only parameter of this section is  $\alpha$ . It can be set as constant through all the training or it can be learned thanks to the approach described in [?]. Further details in section 2.2.5.

**Replay buffer** See section 2.2.4 on page 29.

**Neural Network** See section 2.2.4 on page 29.

**Learning Update** See section 2.2.4 on page 29.

---

**Algorithm 2.5:** Soft Actor-Critic [?]

---

**Input:** Initial policy parameter  $\theta$  and Q-function parameters  $\phi_1, \phi_2$

- 1 Initialise target network weights  $\bar{\phi}_1 \leftarrow \phi_1, \bar{\phi}_2 \leftarrow \phi_2$
- 2 Initialise an empty replay buffer  $\mathcal{D}$
- 3 **for**  $episode = 1, M$  **do**
- 4     Receive initial state  $s_t \leftarrow s_1$
- 5     **repeat**
- 6         Observe state  $s_t$  and select action  $a_t \sim \pi_\theta(\cdot|s_t)$
- 7         Execute  $a_t$  and obtain tuple  $(r_t, s_{t+1}, d_t)$
- 8         Store  $(s_t, a_t, r_t, s_{t+1}, d_t)$  in replay buffer  $\mathcal{D}$
- 9          $s_t \leftarrow s_{t+1}$
- 10     **until**  $s_t$  is terminal
- 11 **end**
- 12 **if** it is time to update **then**
- 13     Sample random minibatch of  $N$  transitions  $(s_t, a_t, r_t, s_{t+1}, d_t)$  from  $\mathcal{D}$
- 14     Calculate targets
- 15      $\phi_i \leftarrow \phi_i - \lambda_Q \nabla_{\phi_i} J_Q(\phi_i)$ , for  $i \in \{1, 2\}$
- 16      $\theta \leftarrow \theta - \lambda_\pi \nabla_\theta J_\pi(\theta)$
- 17      $\alpha \leftarrow \alpha - \lambda \nabla_\alpha J(\alpha)$
- 18      $\bar{\phi}_i \leftarrow \tau \phi_i + (1 - \tau) \bar{\phi}_i$ , for  $i \in \{1, 2\}$
- 19 **end**

**Output:** Optimised policy parameter  $\theta$  and Q-function parameters  $\phi_1, \phi_2$

---

## 2.3 Related Work

A truly inspiring work for this thesis is [?] where the authors show, probably for the first time, that deep reinforcement learning is a viable approach to autonomous driving. Nowadays, most approaches focus on formal logic which determines driving behaviour using annotated 3D geometric maps: the external mapping infrastructure intuitively makes this approach limited to the models and the representation of the surroundings. This technique is not able to scale up efficiently because of last-mentioned strong dependencies.

The fundamental concept underlying [?] to make autonomous driving systems a ubiquitous technology is the design of a system which can drive relying - just like humans - on a comprehensive understanding of the immediate environment [?]. [?] is useful to motivate the research in this direction because it represents an

example of autonomous vehicle navigation exploiting GPS for coarse localisation and LIDAR to understand the local scene instead of detailed prior maps. Many sensors have been developed through the years to gather information and observations increasingly sophisticated. However, the major problem is the massive budget needed to afford all these technologies. The extraordinary results obtained by [?] are not based on intelligent sensing techniques, but on the usage of a monocular camera image together with vehicle speed and steering angle. They decided to apply the model-free approach of reinforcement learning because it is exceptionally general and useful to solve tasks that are complex to model correctly. As discussed previously in this chapter, model-based algorithms tend to be more data-efficient than model-free ones; however, the quality of the adopted model limits the results [?].

The authors decided to exploit Deep Deterministic Policy Gradient (DDPG) [?]. Firstly, they developed a 3D driving simulator using Unreal Engine 4 to tune reinforcement learning hyperparameters such as learning rates and number of gradient steps to take after each training episode. Then, they tried to apply the DDPG algorithm in the real world using the parameters learnt. They also did some experiments using a compressed state representation provided by a Variational Autoencoder [?, ?]. The excellent results obtained in [?] have been exceeded by the same authors in [?] that shows astonishing performances driving on narrow and crowded urban road never-seen during training.

Starting from all these outlined ideas, we decided to investigate ways and approach to autonomous driving with reinforcement learning without using simulators or prior data in order to make a step towards the so-called *reinforcement learning in the wild* [?]. The most popular and prominent achievements in reinforcement learning to date consist of experiments done using simulated environments or ones which exploit the knowledge acquired in simulated environments in real ones. The team of the DeepMind manages to develop smart agents capable of performing super-human results in numerous games and videogames such as *Go* [?, ?], while OpenAI engineers an agent capable of beating the world champion of the multi-player Dota 2 game [?, ?]. These problems are not easy to solve but have the benefit of having a training environment equal to the test one. This fact does not apply to the approach which exploits simulated experiments results in real environments.

The last-mentioned method encloses a critical caveat: the awareness that the simulator has limits. Unlike what happens in some environments where there are many similarities between the simulator and the environment – e.g. games and videogames –, in harder environments – e.g. autonomous vehicle on public roads, ads marketplaces, biology, and applications around human behaviour – is very difficult to get a simulator capable of reproducing with high-fidelity the wild environment because of approximations that could mislead the learning. On the other hand, learning directly on the environment need a fast learning cycle because of the significant number of interaction with the environment and, above all, cheap or

low-risk exploration cost: millions of self-driving episodes in simulator ending with a crash of the car has a small cost in respect of one episode with the same ending crash in the real world.

Between the end of 2018 and the beginning of 2019, UC Berkeley and Google developed jointly a state-of-the-art off-policy model-free reinforcement learning algorithm called soft actor-critic (SAC) (see section 2.2.5 on page 30). The critical goal is to provide a deep RL algorithm suitable for real-world problems and the related new challenges that they pose. They outlined the desired properties of the ideal deep reinforcement learning algorithm:

**Sample Efficiency:** the process of learning in the real world can be a very time-consuming task. For this reason, it is desirable a functional sample complexity to learn skills successfully.

**No sensitive hyperparameters:** as mentioned before, hyperparameter tuning could be a tough task to complete in real-world experiments as it could require numerous repetitions in situations where the cost in terms of time and money could be burdensome. The proposed solution of the authors is maximum entropy RL which provides a robust framework that minimises the need for hyperparameter tuning.

**Off-policy learning:** this learning approach allows the reuse of collected data for various tasks. This possibility helps the process of prototyping a new task when adjusting parameters or choosing a different reward function.

Soft actor-critic (SAC) represents a deep RL algorithm capable of satisfying the previously mentioned requirements. The authors of [?, ?] revealed that the algorithm is capable of solving real-world robotic tasks in a conspicuous but acceptable number of hours, showing its robustness to hyperparameters and working on a large variety of simulated environments always using the same set of parameters. Not only it shows great performances in numerous challenging tasks compared to deep deterministic policy gradient (DDPG), twin delayed deep deterministic policy gradient (TD3) [?] and proximal policy optimisation (PPO) [?], but reveals its power by solving three tasks from scratch without relying on simulators or demonstrations [?].

The real-world robotic task involving a 3-finger dexterous robotic hand to manipulate an object similar to a sink faucet with a coloured end presented in [?, Section 7.3] shows similar concepts to the task analysed in this thesis for what concerns the framework of deep reinforcement learning exploited. The algorithm exploits raw RGB images and processes them through a convolutional neural network. The final goal of the robot is to rotate the valve into the correct position – with the coloured part pointing to the right – starting from a random position for each episode. The results obtained represents one of the most sophisticated



real-world robotic manipulation tasks learned end-to-end with deep reinforcement learning, starting from raw images and without any previous simulation or pre-training phase.

Taking all arguments into account, we decided to follow [?] implementing a similar self-driving framework based on the design of a control system for a small toy robot called Anki Cozmo (see chapter 4 on page 54). Therefore, we formalise the autonomous driving learning problem as a Markov decision process to enable the application of reinforcement learning algorithms, taking into account the improvements mentioned in this section about model-free reinforcement learning algorithms applied to real-world robotic problems. Chapter 4 on page 54 provides a detailed description of the choice made and the system design.

# Chapter 3

## Tools and Frameworks

This chapter aims to describe the main tools and frameworks used to develop the project of this thesis. The first section will describe OpenAI Gym framework, a central toolkit for developing and comparing reinforcement learning algorithms, and explain why this tool is essential for reinforcement learning research. The second part of this chapter will outline Cozmo, the powerful toy robot developed by Anki that we used as an agent in the reinforcement learning scenario to apply algorithms to try solving autonomous self-driving tasks: chapter 5 on page 60 will report the analysis and description of these experiments. This section will also report a set of available alternatives to Cozmo, explaining the motivations underlying the final choice. The description about PyTorch, an optimised tensor library for deep learning using GPUs and CPUs used to build up the convolutional neural network of this work, and the related comparison with TensorFlow will occupy the last part of this chapter.

### 3.1 OpenAI Gym

Nowadays, OpenAI Gym, released in 2016 with its public beta, is one of the most popular toolkits and frameworks in the reinforcement learning scenario. A brief analysis of reinforcement learning research could be useful to outline the motivations underlying the need for a reinforcement learning framework.

As reported previously in chapter 2 on page 2, reinforcement learning is a sub-field of machine learning dedicated to the world of decision making and motor control: researchers study how an agent can learn and improve to achieve a specific goal in a complex, usually unknown environment. This machine learning paradigm is becoming more and more attractive for both researchers and industries because of its visionary property of being very general. A reinforcement learning algorithm can be exploited to control a robot's motor in order to make it capable of running or jumping, play a videogame or a board game, make critical business decisions like

pricing and inventory management, but also learn how to invest in financial trading environments. The generality of reinforcement learning became engaging thanks to the remarkable results achieved in many challenging environments, as reported previously in chapter 2 on page 2.

Despite these appealing features, the research was slowed down by other circumstances, no less critical. The need for better benchmarks represents the first factor. As an example, the abundant availability of conspicuous datasets like *ImageNet* [?] has driven supervised learning improvement in the research. For what concerns reinforcement learning, the nearest equivalent to supervised learning datasets would be a broad collection of different environments in order to test various algorithms with different kinds of observations or rewards. The second drawback of this approach to learning is the lack of standardisation of environments designed in publications. In reinforcement learning, subtle differences in problem definition, reward function design or action space typology could make the difficulty of the task grow. This fact threatens to slow down and corrupts experiments reproducibility making an objective comparison between the results of different papers almost impossible.

The need to fix both problems was the primary motivation behind the design and implementation of OpenAI Gym.

### 3.1.1 Environments

The agent and the environment represent the main components of reinforcement learning. The choice of OpenAI was to implement and provide the abstraction mainly for environments, not for agents. They decide to provide a standard environment interface instead of forcing the developer to use pre-defined agent interfaces: the motivation behind this choice was to leave developers independent in the design of the agent, the core of reinforcement learning, and facilitate the creation and usage of environments. Thanks to this approach, all agents implemented with OpenAI Gym can be used with the whole set of environments provided by the framework. Therefore, it is possible to create a personalised environment to suit the needs of a specific experiment that can be used by all agents exploiting OpenAI Gym environment interfaces.

In this scenario, we realised the first contribution to our thesis. Thanks to this framework features, we implemented an OpenAI Gym environment capable of interacting with Anki Cozmo by providing a binding between functions of Cozmo SDK and interfaces of the reinforcement learning framework. In chapter 4 on page 54 we will provide further information and details about our contribution.

The importance related to the high quantity of environment is fundamental to build a reliable and sustainable framework for reinforcement learning algorithms. For this reason, OpenAI Gym contains a various and heterogeneous environment database, ready to be used.

## Interface Functions

Exploring OpenAI Gym, it is essential to focus on the most crucial interface functions that the agent will exploit to interact with the environment. The functions which constitute the skeleton of an OpenAI Gym environment are the following:

- **def step(self, action):** through this function, the agent can communicate the action it wants to take. The input data depends on the type and number of variables in the actions space (e.g. discrete or continuous). As will be discussed in section 3.1.2 on the next page, the values returned by this function represent the environment state after the manipulation caused by the agent action. Thanks to these data, the agent will be able to select the next action following the reinforcement learning loop.
- **def reset(self):** during the episode, internal variables of the environment changes, influenced by the action taken previously. This function allows the agent to restart the initial situation of the environment. This procedure is particularly helpful when an episode finishes and the agent has to restart the next learning episode in a brand new copy of the environment.
- **def render(self, mode='human', close=False):** this function is mainly used in simulated environments. It enables the visual render (if available) of the environment.
- **def close(self):** the final function to close the environment after the end of all experiments and episodes.

## Available environments

To date, OpenAI Gym includes the following environments:

- **Algorithms:** learning to imitate computations, such as copying or reversing symbols from the input tape, is the main aim of this typology. These environments might seem easy to be solved by a computer, but it is important to remember that the objective here is to learn to solve these tasks purely from examples. Therefore, it is possible to vary the sequence length to increase or decrease task difficulties easily.
- **Atari:** *Atari 2600* is a home video game console developed in 1977 which spread the use of general-purpose CPUs into gaming with game code distributed through cartridges. This environment section provides a database which contains more than 100 environments emulating Atari 2600 videogames. OpenAI Gym exploits *Arcade Learning Environment (ALE)* [?] providing RAW pixel images or RAM as observation of the environment.

- **Box2D**: in this group is possible to find some continuous control tasks in a simple 2D simulator such as *BipedalWalker*, *CarRacing* and *LunarLander*
- **Classic Control**: this class provides a set of problem borrowed by control theory and widely exploited in the classic reinforcement learning literature. Some task examples are balancing a pole on a cart or swing up a pendulum.
- **MuJoCo**: this collection contains continuous control tasks running in a fast physics three-dimensional simulator called *MuJoCo* which stands for *Multi-Joint dynamics with Contact*. This physics engine aims to facilitate research and development in robotics, biomechanics, graphics and animation. The simulator is particularly suitable for model-based optimisation allowing to scale up computationally-intensive techniques. Thanks to its features, it became useful as a source for reinforcement learning algorithms. [?]
- **Robotics**: OpenAI released this algorithm typology to provide eight robotics environments with manipulation tasks significantly more difficult than the MoJoCo ones. It contains *Fetch*, a robotic arm to move and grab objects, and *ShadowHand*, a robotic hand to manipulate and grab pens, cubes and balls. [?]

### 3.1.2 Observations

As previously reported, the `step(self, action)` environment interface is the most important one because it contains the behaviour definition of the environment that reacts to agent actions. Indeed, the agent has to know in which way its actions are influencing the environment in order to stop doing random actions and start making valuable decisions.

In order to provide this type of information to the agents, the `step` function returns four relevant values that reinforcement learning algorithms can exploit to determine the best action to do in the future. These values are:

- **observation (object)**: a specific object which represents the environment observation and shows the changes provoked by agent actions. Its structure and interpretation depend on the implementation of the specific environment. For example, it could represent the raw pixel data from a camera, the status of a board game or physics data of a robot (joint angles and velocities).
- **reward (float)**: this is the fuel of the reinforcement learning algorithms. This value represents the reward achieved thanks to the action taken by the agent. The reward for each action can change among environments, but this is the crucial information to enable the agent to learn.

- **done (boolean)**: this value is a simple boolean that signals the agent when the episode ended and it is time to reset the environment to start a brand new episode.
- **info (dict)**: a Python dictionary to monitor and show diagnostic information useful for debugging. The agent could use it as additional information in the learning process, but the documentation of OpenAI Gym does not recommend to use this data in the learning process to maintain the coherence with the reinforcement learning loop.

The last important thing to report about OpenAI Gym framework is the definition of **Spaces**. As reported in the second chapter, an environment consists of the action space and the observation space. Both can provide discrete or continuous values: this distinction is fundamental to determine the usage of a specific algorithm rather than others to solve a given task better. For this reason, OpenAI Gym provides **Discrete** and **Box** to initialise discrete and continuous spaces, respectively. It allows the user to define a lot of features useful in the learning process, such as a maximum and minimum value setup or a default function to instantly sample a random value from the defined space.

## 3.2 Anki Cozmo

Especially in the latest decades, human beings have started a complicated relationship with robots. They are very fascinated by the prospect of artificial intelligence offered by recent researches and applications. However, they are apprehensive and worried at the same time because of the apocalyptic plot that many sci-fi films show and the big promise of automation, capable of replacing human workers in the future. However, all these concerns disappear after meeting Cozmo, the palm-sized toy robot developed by the San Francisco-based company Anki and available on the market since 2016. On first glance, Cozmo might appear as one of the cutest toy robots: not surprisingly Anki employed the guidance of Carlos Baena, the former Pixar animator, to design this robot. Therefore, it can interact with people using a small display screen together with audio effects to mimic human emotional reactions and responses. The result is a toy robot WALL-E-inspired both aesthetically and personality-wise, powered up by artificial intelligence to move and discover the surrounding environment. Thanks to the built-in camera, Cozmo can remember faces and recite names, but also to plan paths and play various games with its three cubes that carry sensors and lighting.

Despite these entertaining, but not so technical facts, Cozmo hides a lot of powerful features under the hood. Anki developers produced a high-quality Python SDK that allows developers to take control of the whole set of Cozmo sensors and actuators thanks to interaction granularity offered by functions and interfaces.

This section aims to outline the hardware and software architecture hidden underneath the cute Cozmo bodyworks, together with a comparison with the alternatives available to design the reinforcement learning system of this thesis. An essential inspiration in the writing of this section came from [?], [?] and Cozmo SDK Forums <sup>1</sup>.

### 3.2.1 Cozmo Architecture

It is possible to define Cozmo as a vision-guided mobile manipulator, one of the first consumer robot which can boast vision among its features. However, as reported before, the features it shows during the normal toy usage, do not equal the number of interfaces and functions made available to developers. The hardware and the software developed for Cozmo makes it the right choice to fast prototyping computer science projects: it is the main reason why we decided to exploit Anki Cozmo instead of other alternatives. The SDK provided by Anki consists of a comprehensive set of low- and high-level functions which grants full access to sensor data providing the right flexibility, simplicity and granularity to satisfy every developer needs. It is versatile because it could be easily connected with hundreds of third-party libraries to augment Cozmo capabilities. Therefore, it is an entirely open-source SDK to give the community the freedom to customise and contribute.

Figures 3.1 and 3.2 on the following page and on page 46 shows technologies, the hardware and software involved in the production of Cozmo. The first image represents stacks and connections between the robot core and the mobile application provided by Anki, while the second one shows the interaction between the last-mentioned application and the Python user program on the personal computer. The reader can retrieve a global perspective about the whole Cozmo architecture by merging these two figures.

Starting from fig. 3.1 on the following page, it is noticeable that Cozmo has a lot of sensors and actuators, which enable it to move and understand the surrounding environment. For what concerns the sensor part, the VGA Camera is the crucial component exploited in the thesis. This camera provides a grayscale image with  $320 \times 240$ . Anki reported a camera resolution equal to  $640 \times 480$ , but the latest firmware version supports only the lower resolution. Furthermore, the camera can detect and acquire colours, but the firmware limits this feature to maintain the bandwidth stable and avoid overloads or slowdowns in the communication with the mobile application and subsequently, the user program running on the computer. The camera has approximately  $60^\circ$  field of view (FOV) and 290mm focal length.

As regards the actuators column, Cozmo can explore the world thanks to its four motors and over fifty gears. Instead of having wheels, this robot has two tracks

---

<sup>1</sup>Link: <https://forums.anki.com/>

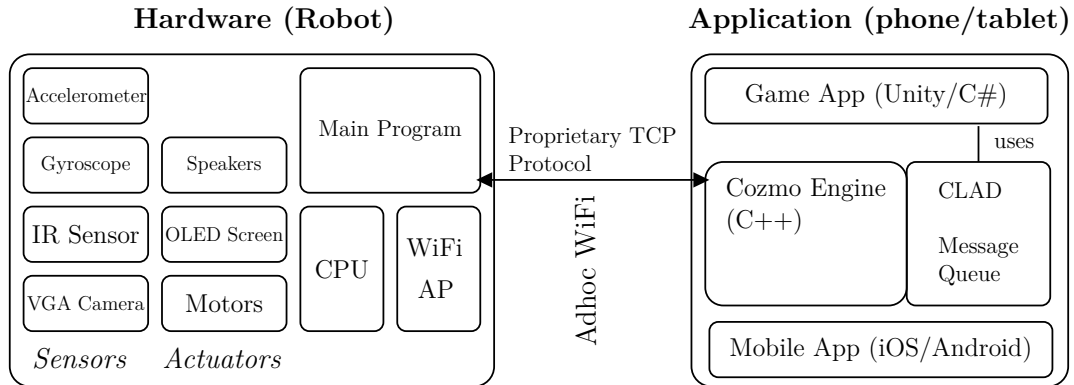


Figure 3.1. Interaction between the Robot and the mobile application stack. The robot main program interacts with the Cozmo Engine (C++) implemented in the mobile application through a proprietary TCP protocol established using and ad-hoc WiFi connection. [?]

to navigate: it can steer and move freely by controlling the speed of each track. Another moving part is the head that can move up or down to direct the camera and the display. It also has a forklift with which it can lift objects or the cubes available in the kit.

Cozmo can perform one single action at a time: if the current action is not yet complete and the request of a new action occurs, the new action fails with a *tracks locked* failure code. For this reason, the SDK does not allow the execution of multiple actions. Despite this fact, it is easy to notice that Cozmo animations and behaviours typically use combinations of actions: indeed the user can send simultaneous actions to the robot by calling them using `in_parallel=True` in its script, but these actions have to belong to a different action track. This approach allows the parallel execution of actions, provided that they use different tracks. Cozmo software architecture holds seven independent action tracks:

- **HEAD**: raise/lower robot head.
- **LIFT**: raise/lower robot forklift.
- **BODY**: wheels or treads actions for driving and turning.
- **FACE\_IMAGE**: actions with the OLED display such as animations or faces.
- **EVENT**: for this action type Anki does not release further information.
- **BACKPACK\_LIGHTS**: actions or animations of lights on Cozmo back.
- **AUDIO**: speech or sound effects emitted by the robot.



The project of this thesis mainly exploited the body action track to move the robot in the environment: it also employed head and lift tracks, but only to position the robot head to provide images about the track to the main program and the reinforcement learning agent. The usage of parallel operations was not necessary.

The Cozmo hardware includes an onboard CPU and a WiFi access point, thanks to which the user can interact with the robot. However, to activate this communication, it is necessary to install the Cozmo Android/iOS application on a personal tablet or smartphone. This application is the same used to play with the toy-side of Cozmo, but in its settings, there is a function to enable Cozmo development mode. After connecting the chosen personal device to the robot through a simple WiFi connection, the application manages to create a proprietary TCP protocol to interact directly with the main robot program. The creators of Cozmo developed this mobile application using C++ to implement the *Cozmo Engine*, the component which shall be responsible for the TCP communication with the robot. To design and implement the game experience and the graphical user interface, they used Unity and C#.

This last-mentioned part utilises the component internally designed and implemented by Anki to provide a contact point between the Cozmo SDK installed in the development machine, as shown in fig. 3.2 on the next page: the C-like Abstract Data language (CLAD). The fundamental idea behind this tool is to make the process of serialisation, communication and deserialisation of data structures written in Python easier for the developer. In practice, for every data which has to pass over the wire, files with extension `.clad` to define enums, structures and messages are generated: their syntax is similar to C struct one. After that phase, this tool auto-generates Python, C++ and C# code for each structure previously defined. This process allows the user to define a specific message in Python and to send it over the network to C++ where it will be deserialised automatically, avoiding problems and sources of bugs coming from intricate underlying details.

This approach results in a process method much lighter than the one provided by *Protocol Buffers (Protobuf)* by Google: even then, the main aim of this implementation choice was reducing the network bandwidth. Taking all arguments into account, CLAD is the protocol used for the interaction between the SDK and the C++ Engine that allows the developer to not care about the low-level part of the code and to focus on the high-level logic using Python. It is also useful to maintain the same interface for the user even if low-level logic changes occur.

In this case, the CLAD communication messages exploit a wired connection between the mobile device and the computer with a simple USB cable. In this thesis, we used an Android tablet to run the experiments, and for this reason, we used the *Android Debug Bridge (ADB)* to make this exchange possible. It is a command-line tool included in the *Android SDK Platform-Tools* package that allows the communication between the computer and Android devices. It facilitates many device actions such as installing, debugging apps or running a variety of

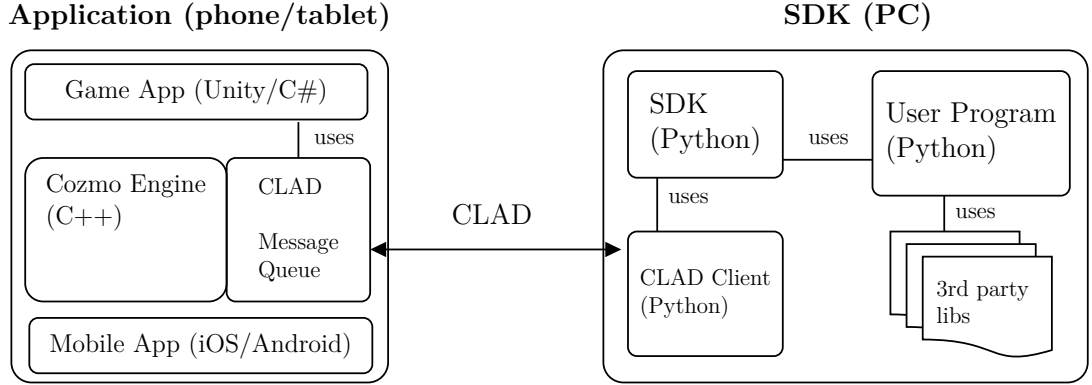


Figure 3.2. Interaction between the mobile application stack and the Python user program. Anki designed a C-like Abstract Data Language (CLAD) to implement the connection between the Cozmo SDK and the Cozmo Engine on the mobile application. This approach puts a decoupling layer between Python interfaces and low-level functions to make prototyping easier and fast-forwarding for developers. [?]

commands on a device. It is a client-server program that includes three components: it includes a *client* hosted in the development machine, a *daemon* (*adbd*) which runs the commands on the device as a background process and finally a *server*, a background process in the development machine that manages the communication between the two previous components.

The last step to start programming with Cozmo SDK is to install it in a development machine following the instruction provided by its documentation <sup>2</sup>.

### 3.2.2 Why Cozmo?

Before starting the development of this thesis project, we spent some time analysing the small car market situation in order to find out the right choice for our specific needs. In addition to Anki Cozmo, the ideal alternatives in the self-driving scenario when we started this thesis were *AWS DeepRacer Vehicle* and a customised car implemented from scratch exploiting *Donkey @Car library*.

This section part aims to briefly describe the main alternatives to Anki Cozmo listing strengths and weaknesses for each approach to motivate the final choice made.

<sup>2</sup>Cozmo SDK Documentation: <http://cozmosdk.anki.com/docs/>

## AWS DeepRacer

*AWS DeepRacer* is a platform developed by Amazon to learn and refine machine learning and reinforcement learning algorithms and techniques on AWS DeepRacer vehicle.

The whole ecosystem consists of two main part:

- **AWS DeepRacer Car.** It is a 1/18 scale race car developed and built to test reinforcement learning algorithms on a real track. It aims to show how the reinforcement learning model trained in a simulated environment can be conveyed to the real-world by using cameras to view the track and the model to control throttle and steering. It has a lot of interesting specifications under the hood such as Intel Atom<sup>®</sup> Processor, 4GB of RAM, an expandable 32GB of storage to accommodate the trained model and a 4-megapixel camera with MJPEG.
- **AWS DeepRacer Simulator.** The user can build his models in *Amazon SageMaker* and train, test and iterate the learning process using the racing simulator. It offers an integrated environment hosted on the AWS cloud to experiment and optimise algorithms to apply to the autonomous driving model.

There are many advantages in using the DeepRacer stack. It offers an integrated approach where the developer has to focus just on reinforcement learning: indeed, it abstracts a significant part of the software giving the developer the chance to work almost exclusively on the model. It has a growing development community to confront ideas and find solutions through challenging competitions around the world. One of the fundamental aims behind this product is to build up a regulated environment where amateurs and researchers can compare and tests different approaches in solving the autonomous driving task. Therefore it provides a high-performance car where to store the trained model in order to make it independent to the training environment.

On the other hand, one of the main disadvantages of this approach is that it is an Amazon lock-in: the framework provided force developer to use the whole Amazon stack. This fact can be inconvenient economically-wise because of the rates of the Amazon platform that adds to the price of the car, but, above all, because Amazon may have access to all developer implementations and experiments.

Another not negligible factor was the product release date: Amazon had to release it in July, but, to date, it is not yet available in Italy. Therefore, because of this thesis is focused on the application of reinforcement learning algorithms directly in the real world without the aid of simulators and model, the absence of the physical car was crucial for the final decision.

Another important withdraw consists of the strict correlation between the simulator and the real system: the framework provided by Amazon seems to be suitable

only to test reinforcement learning algorithms after the model training in the simulator, but not for the sake of this thesis.

### **Donkey®Car**

The second important alternative to Cozmo was to develop a small car from scratch using integrated boards such as *RaspberryPi* or *NVIDIA Jetson Nano*, sensors and camera to personalise the device and satisfy specific project needs. *Donkey®Car* <sup>3</sup> is the most popular choice to build a personal self driving toy car. It is an open-source community project powered by volunteers who are interested in build their own self-driving cars. They cooperated to build up a high-level self-driving library written in Python, focusing on enabling fast experimentation and easy contribution. They provide detailed instructions to build a personal Donkey, providing kits or lists of components to use: the whole kit costs among 250-300\$.

It is also possible to install the Donkey library on any RC car to make it self-driving and autonomous. Despite this fact, Donkey developers suggest building the *Donkey2* car, which is a tested hardware and software setup, to avoid problems such as incompatibilities or bugs and make the most out of the presence of a dense community.

The main strength of this choice is the freedom to develop a completely customised car to better suit the needs and requirements of a great set of autonomous driving projects. In the last releases of the library, Donkey developers released a complete sandbox simulator for training a self-driving car. The languages used to develop it are Unity for simulation and Python with Keras and Tensorflow for training. They also provide an OpenAI Gym environment to use with such simulator.

The main weakness of this method consists in the process of building a self-driving car system from scratch. It offers versatility and flexibility in components choice, but the time to devote to building a working system, free from as many bugs as possible would have slowed down the prototyping of the system itself and the whole thesis project. The duration of development is not associated with the process of physically assembling the car because Donkey developers estimate two hours to build the car. The real obstacle consists of the setup of a connection similar to the Cozmo one: it needs to be as stable as possible to make all stack working smoothly in real-time.

### **The final choice**

As easily predictable, the final choice fell on Anki Cozmo. It provides a fast-forwarding SDK ready to be exploited in prototyping a brand new project together

---

<sup>3</sup>Donkey Car community website: <https://www.donkeycar.com>

with the strictly necessary sensors to perform the reinforcement learning experiments in the real world. A crucial factor in reaching the final decision was the dimension of the car: as reported before, both Donkey car and Amazon DeepRacer are 1/10 and 1/18 scale race car respectively, while Cozmo is just 5.5cm wide, which results in a ratio of about 1/30. This fact not only makes Cozmo an easily transportable solution to speed up experiments and to restart episodes effortlessly but also a solution that supports the design of a track in a restrict space such as the one in the laboratory of Eurecom.

Therefore, the connection between the development machine and the robot is suitable for implementing an OpenAI Gym environment, and it is similar, at least in the premises, to the distributed off-board computation approach. The main algorithm, the neural network, the reinforcement learning framework and the other cognitive parts are computed and managed by the workstation, instead of being stored in the vehicle as happens in Amazon DeepRacer and Donkey car approaches.

Taking into account the perspective of autonomous cars in the real world, on-board and off-board computation approaches are still under research. With the on-board method, cars have much computational hardware inside in order to manage every aspect of autonomous driving by themselves, but it requires more powerful batteries to counterbalance energy consumption. On the other hand, the connection to off-board computer facilities or the clouds leads to new vectors of attack but also enables companies to monitor the behaviour the vehicle fleet to identify malicious activities early. To date, both approaches are still under research, and there it is not possible to decree a legitimate winner. For what concerns the thesis, the designed system emulate an off-board approach.

In the end, Cozmo provides plain and straightforward control of the car and a rich Python SDK to use with OpenAI Gym and it is the best trade-off between functionalities and fast-developing.

### **3.3 PyTorch**

PyTorch <sup>4</sup> [?] is an open-source machine learning and deep learning library developed by Facebook's AI Research Lab and released to the public in October 2016. The main aim of PyTorch is to provide an intuitive and straightforward framework to develop artificial intelligence projects: two of the main applications to date are computer vision and natural language processing.

The programming languages utilised to develop PyTorch were Python, C++ and CUDA, the parallel computing and API model created by Nvidia to allow software developers and engineers to use CUDA-enabled GPU for general purpose processing. The primary interface provided by the library to the user employs

---

<sup>4</sup>PyTorch Github Repository: <https://github.com/pytorch/pytorch>

Python, the project where Facebook developers mainly put their efforts. Despite this fact, it also offers a C++ interface.

PyTorch consists of the following components:

- **torch**: PyTorch Tensor library with strong GPU support. It implements interfaces similar to those of the NumPy library. It contains data structures for multi-dimensional tensors and mathematical operations, providing many utilities for efficient serialising tensors and arbitrary types.
- **torch.autograd**: the tape-based automatic differentiation library that supports every differentiable operation on tensors available in **torch**.
- **torch.jit**: this component is a compilation stack that uses TorchScript to create serializable and optimizable models from PyTorch code. This tool allows the user to train models in PyTorch using Python and then export the model in a production environment where Python may be disadvantageous for performance and multi-threading reasons.
- **torch.nn**: this component provides a neural networks library that is entirely compatible with **autograd** and designed for flexibility.
- **torch.multiprocessing**: this component is based on the Python multiprocessing library, but it implements memory sharing of torch tensors across processes.
- **torch.utils**: it contains many utility functions to better exploits the features of PyTorch.

PyTorch provides a NumPy-like experience to interact and manipulate data structures suitable for GPU computation, offering a deep learning research platform which can provide flexibility and speed. These data structures are called *Tensors* and they can be used both on the CPU and the GPU, accelerating the computation thanks to the whole set of functions and manipulators explicitly designed for every scientific computation need.

PyTorch is not a straightforward binding to an underlying complex C++ framework. The library design focused on establishing Python as the main priority, and for this reason, the user experience is very natural and similar to other important machine learning libraries already present in the package manager.

It is noticeable that PyTorch developers aimed to create an intuitive and linear product to use. To follow this idea, they decided to make PyTorch synchronous to permit the debugger to receive and understand messages and stack traces promptly. This feature translates in a better debugging experience for the end-user.

Beyond these features, one of the traits that distinguish PyTorch from other frameworks is its single way to build neural networks by using a tape-based automatic differentiation. The majority of deep learning frameworks available in the

market, such as TensorFlow, Theano or Caffe, exploits a static approach in structure computation graph creation: they reuse the same layout in the whole program, therefore changing a simple component triggers the regeneration of the graph from scratch. PyTorch utilises an entirely different approach which is not unique to PyTorch, but it provides one of the fastest implementations: they call it *Tape-Based Autograd*. This term refers to the reverse-mode automatic differentiation exploited in the framework, which is a technique based on the properties of the chain rule: to calculate the derivative of an output variable w.r.t. any intermediate or input variable, the only requirement is to know the derivatives of its parents and the formula to calculate derivative of primitive expression. The main improvement that this approach brings is allowing the user to change the network structure on-the-fly without lag or overhead.

### 3.3.1 TensorboardX

One of the most crucial means that every machine learning researchers need is a tool to visualise and measure data efficiently: this fact is significant because to improve models, projects and results, we need to measure. However, one of the main problem in PyTorch is the absence of such a tool, specifically designed for the Facebook framework. It can always use powerful tools such as *Matplotlib*, but it offers a synchronous approach that leads to slow down the main program since the primary process has to wait for the data rendering before starting next operations.

On the other hand, TensorFlow, the most important deep learning alternative to PyTorch developed by Google, provides in its package TensorBoard which is a webserver to serve visualisation of the training progress of the neural network. It can show to the user scalar values, images or text, and it is particularly useful to visualise experimental metrics such as loss and accuracy. The particularity of this tool consists in the fact that it stores these typologies of information asynchronously as events. The Python script calls specific functions to store information and goes on with the next operation without waiting for its render: that is possible thanks to the decoupling level inserted by this approach between the visualisation and the creation of data. Indeed, Tensorboard operates by opening and reading TensorFlow events files that contain the summary data generated during experiments. Therefore, it will be the webserver to take care of elaborate data for rendering without bothering the current Python script.

Fortunately, PyTorch can exploit the features of Tensorboard thanks to a library called *TensorboardX* <sup>5</sup> that stands for *Tensorboard for X* to highlight developers aim to make Tensorboard available for all deep learning framework.

---

<sup>5</sup>TensorBoardX documentation: <https://tensorboardx.readthedocs.io/>

### 3.3.2 PyTorch vs. TensorFlow

After the advent of deep learning, many companies decided to put their efforts to design architectures and frameworks to vehiculate this new technology. The two most popular frameworks in this research field are TensorFlow [?] by Google released in 2015 and PyTorch [?] by Facebook released in 2017.

Implementing the same neural network in these two frameworks will lead to different results because of the training process has many parameters that depend on the underlying technologies provided by the specific framework. For instance, the training process in PyTorch is enhanced by CUDA GPU usage, while TensorFlow can access to GPU through its GPU acceleration. The choice between these two frameworks is not straightforward because it depends on the perspective and the needs of the specific projects to develop. For this reason, this section aims to outline differences between these two libraries without aiming to decree the best one but to motivate the decision to use PyTorch.

#### Dynamic versus Static

The first difference concerns the construction of the computational graph. A computational graph is an abstraction useful to represent the computation process through a direct graph.

In TensorFlow, computational graphs are defined statically, before running the code. The main advantage of this method is allowing parallelism and dependency driving scheduling, features that boost the learning and make it more efficient. This framework communicates with the external world via specific tensors that will be substituted by input data at runtime. Only with TensorFlow 2.0, Google decided to implement dynamic computational graph in its product, but its stable version was released after the start of this thesis.

As mentioned above, PyTorch approach to computational graphs is dynamic. This characteristic means that the graph is built incrementally at runtime without using particular data structures as placeholders. This feature supports projects where the author needs to change the computational graph on-the-fly avoiding the application restart. In this sense, PyTorch is more pythonic than TensorFlow.

#### Distributed Training

Another key feature is the distributed training and data parallelism. PyTorch offers native support for asynchronous execution from Python, and then it could improve performances. On the other hand, TensorFlow needs more efforts to allow distributed training: the developer must fine-tune every computation to make it running on a specific device. Both frameworks offer the same opportunities in these terms. However, TensorFlow needs more effort to make things work.



## Visualisation

As discussed in the previous section, TensorFlow exploits TensorBoard to provide all tools that machine learning researcher needs to visualise learning and keep track of the training process. Facebook researchers developed *Visdom* for this purpose, but it provides very minimalistic and limited features compared to the ones offered by TensorBoard. As reported before, it is possible to use TensorBoard with PyTorch thanks to the library TensorBoardX.

## Production Deployment

For what concerns the deployment of trained models into production, TensorFlow offers the best service via *TensorFlow serving*, a framework that offers and uses REST Client API. The production deployment in PyTorch improved from its early releases, but, currently, it does not provide a framework to deploy the trained models on the web: the developers must use Flask or Django as backend server to provide the right environment to exploit the model.

## Conclusions

Considering all the points explained in this section, we decided to utilise PyTorch for this project, but it is noticeable that there is no winner in this comparison.

Both frameworks have strengths and weaknesses that depends on the specific applications where we would use them. TensorFlow is a mature and robust tool, notably suggested for production and AI-related products. Although it needs some time to get the developer used to its programming approach and, at least at the start of this thesis project, it supports only static computational graph methods. On the other side, PyTorch is an efficient and young framework with a large community and which provides dynamic computational graphs and is more Python friendly. Therefore, it is especially recommended for research-oriented developers.

# Chapter 4

## Design of the control system

In the previous chapters, we outlined deep reinforcement learning fundamentals and the most critical underlying concepts, and then we discussed the choice made about the technologies to use as baselines for our experiments. The decision fell on Anki Cozmo because of the high-quality SDK provided to developers and the versatility and flexibility provided by PyTorch, especially in a research context.

The next step in this thesis is the merge of reinforcement learning theory with the tool presented previously. Indeed, this chapter aims to describe this merge process that results in the design of the control system for reinforcement learning experiments with Anki Cozmo. The work presented in this section represents one of the contributions of our thesis and the necessary step to start reinforcement learning experiments.

The outline of the whole ecosystem with the description of interfaces, frameworks and technologies used occupies the first section of the chapter. This part also comprehends a discussion about DDPG [?] and SAC [?, ?] implementation with references to the choice made in terms of hyper-parameters and problems faced.

The second part of the chapter aims to describe the implementation of Anki Cozmo OpenAI Gym environment from the problem formalisation as MDP to the implementation of human-robot interaction.

In the final section of this chapter, the design and setup of the real track will be discussed together with a discussion about the problems faced and the choice made to overcome them.

### 4.1 Outline of the system

The development of the control system for Anki Cozmo was the main contribution of the thesis, together with the implementation of DDPG and SAC algorithm. The main aim of this work was to create an OpenAI Gym environment capable of interacting with a robot in the real world, not just through the usage of a

simulator. OpenAI Gym usually provides plain and straightforward interfaces to interact with simulated environments: we decided to exploit these functions to allow the application of reinforcement learning algorithm directly in the real-world decision of the robot.

The fundamental source of inspiration to develop this control system was [?, ?]. This publication represents, as its authors reported, the first reinforcement learning self-driving experiment where a car learned to drive through the application of a reinforcement learning algorithm, by trial and error. They first trained the model exploiting Deep Deterministic Policy Gradient (DDPG) in a simulator for many epochs and, after this learning process, they managed to transfer the knowledge acquired in the simulator in real-world experiments.

We decided to export and implement these ideas in our project, adapting them to the specificities and particularities of the Cozmo setup. The resulting system can be summed up by fig. 4.1 which provide a schematic overview of every technology employed and the interaction among them. This section aims to describe as clearly as possible all the components of the control system we designed.

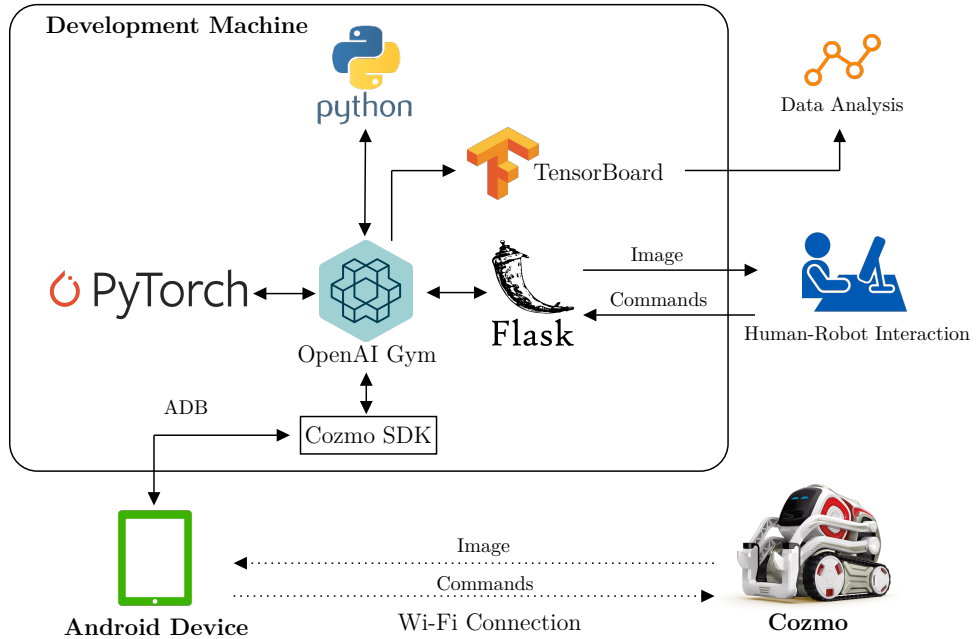


Figure 4.1. Interaction between the user and Anki Cozmo. The main Python script utilises OpenAI Gym for the reinforcement learning component and PyTorch for the deep learning one. The user can influence with the flow of the system through a simple web app that interacts directly with OpenAI Gym and the Cozmo SDK to provide information for the user (e.g. images, learning information) and the robot (e.g. commands). The last component consists in TensorBoard thanks to the script can store results that can be retrieved later by the user.

### 4.1.1 Human-Robot Interaction

It is noticeable that the simulator can undoubtedly be programmed to understand when the car is outside the track. In a reinforcement learning scenario, this fact facilitates the restarting procedure for an episode: the developer has to bind some events or actions to a process that stops the current episode, put the car on the road again and starts the next experiment. In the real world, the situation is more complicated because there are more variables to take into account: the most critical factor is that the failure of an episode in the simulator has no threats or costs, while real-world experiments failures could lead to high costs and severe damage to the car and equipment. Despite this point, the application of such experiment typology in a real environment instead of a mere simulation represents an exciting challenge that could bring reinforcement learning to the next level. In order to make experiments as safe as possible, the authors of [?, ?] implemented a self-driving system designed explicitly for reinforcement learning experiments where the car driver has the faculty of stopping the car when it is going to run off the road or in a dangerous situation and relocating it in the nearest safe position to start the next learning episode.

We aimed to implement this kind of user-robot interaction also in our setup. The requirements were the same as one described in the paper, but this case could not rely on steering wheel, brakes and accelerator as in the car scenario where they are directly accessible by the user. The system needed an interface to allow the user to manoeuvre the robot just like a car.

For this reason, we firstly implemented a straightforward interface using a web app implemented using plain HTML, CSS and Javascript for the frontend and Flask as backend. The aim of this application is allowing user interaction with Cozmo and Flask represented the right choice to allow the communication between this interface and the OpenAI Gym environment. The Flask backend interacts directly with the Cozmo SDK functions to allow the user to see the live streaming from the camera of Cozmo directly in the web app. Therefore it can receive, convert and forward commands to the robot to SDK. These commands consist of the pressure of both single or combination of buttons: Flask decides the action to trigger programmatically following hard-coded conditions and calculates speeds and accelerations of both right and left treads. The SDK function we utilised to implement the manual Cozmo control is `drive_wheels()`. The parameters of this function are the following:

- **l\_wheel\_speed** (*float*): mandatory parameter that specifies the speed of the left tread (in millimeters per second).
- **r\_wheel\_speed** (*float*): mandatory parameter that specifies the speed of the right tread (in millimeters per second).

- **l\_wheel\_acc** (*float*): optional parameter that specifies the acceleration of the left tread (in millimeters per second squared). The default value is the equal to **l\_wheel\_speed**.
- **r\_wheel\_acc** (*float*): optional parameter that specifies the acceleration of the right tread (in millimeters per second squared). The default value is the equal to **r\_wheel\_speed**.
- **duration** (*float*): it specifies the duration of the driving action of the robot. It calls **stop\_all\_motors()** after this duration has passed. The default value is **None**, in this case, the behaviour of this function is equivalent to the non-async **drive\_wheel\_motors()** one: the wheels will continue to move at that speed until commanded to drive at a new speed.

This method is the same we used to implement the OpenAI Gym environment, allowing the algorithm to interact with the robot just like a human in the driving seat. Just as examples, the user can start and stop the episode using the enter button, forget the previous episode with backspace and relocate the robot on the track using W, A, S and D buttons.

In this context, the human-robot interaction has a crucial role in the learning process because it is the source of all improvements and flaws of the learning process: the robot learns which actions are worthy and which are not, but the human is the one who decides the correctness of each action, carrying its unconscious bias in the algorithm. We will discuss the experiments and their results thoroughly in chapters 5 and 6 on page 60 and on page 61.

### 4.1.2 OpenAI Gym Cozmo Environment

This section describes the main features we implemented in *CozmoEnv*, the OpenAI Gym environment we developed to apply and test reinforcement learning algorithms with Cozmo.

#### MDP formalisation

The first important step to implement a reinforcement learning environment is formalising the problem we want to solve as a Markov Decision Process (MDP). Starting from the definition of MDP that we already reported in eq. (2.4) on page 5, we studied and searched the best configuration to provide an acceptable trade-off between performance and memory consumption, capable of making the system work on the development machine we used for our experiments. In the following paragraphs, we will analyse each MDP aspects paying particular attention to the problems we faced and the motivations behind every decision made to overcome them.

**Observations** The environment observation is the group of information that the reinforcement learning agent can see and obtain from the interaction with the real world. In the self-driving tasks in question, the most reasonable way to obtain knowledge about the surrounding environment is to exploit front camera images provided by Cozmo SDK. However, it is essential to remember that the state must satisfy the Markov property according to which it must be independent of the past given the present. Using a single image to describe the current state may be an oversimplification of the problem. Even a human would not be able to determine the best action to take starting from a single frame. That is because it is difficult to infer or determine the current direction of the car, speed or steering wheel position. Merely adding another image can intuitively improve the human perception of the scene.

To verify this fact before starting the experiments with the application of the reinforcement learning algorithms to Cozmo, we firstly analysed the results of some experiments on the environment of the classic control problem of trying to keep a frictionless pendulum standing up (*Pendulum-v0*), appropriately modified to exploit raw images as observations instead of information about its inclination. We noticed that the usage of a single image as observation leads to more unstable and worse results than the ones obtained by combining two subsequential images to feed the algorithm. For this reason, we decided to use two consecutive raw images provided by Cozmo by using a queue data structure with size two. In the code flow, the developer merely has to push the next image obtained by the robot and use the specific queue function implemented to obtain the concatenation of the last two images, only when necessary.

Therefore, we decided to reduce the original size of the input Cozmo image to  $64 \times 64$ . We made this decision as a trade-off between performance, learning phase duration and space available in the central memory to store the replay memory.

In conclusion, the observation provided to the agent is the concatenation of two images with dimension  $64 \times 64$  pixels.

**Actions** Another essential part that we needed to formalise is how the agent can interact and influence the surrounding environment: its actions. In the car scenario, the two main components of driving are the speed and the position of the steering wheel. To straightforwardly formalise these two parts, we decided to use two simple real values.

We chose to define the desired speed value in a range of 0 to 1, while we opted for a range of -1 to 1 for the steering wheel position. These limits allow the agent to manipulate simple values externally, while the underlying logic of the environment takes care to translate these values to a compatible format for the Cozmo SDK. Indeed, as we reported before in this section, the Cozmo SDK function exploited to manoeuvre the robot needs at least the speed of each tread. Algorithm 4.1 on the following page reports the principal steps of this translation.

---

**Algorithm 4.1:** CozmoEnv actions conversion from virtual to real

---

**Input:** Desired speed  $s_t \in \{x \in \mathbb{R} | 0 \leq x \leq 1\}$   
Steering wheel position  $w_t \in \{x \in \mathbb{R} | -1 \leq x \leq 1\}$   
Maximum forward speed  $s_{\max}^{\text{forward}} = 150\text{mm/s}$   
Maximum turning speed  $s_{\max}^{\text{turning}} = 100\text{mm/s}$   
1 Left tread speed:  $ts_{\text{left}} = s_t \cdot s_{\max}^{\text{forward}} + w_t \cdot s_{\max}^{\text{turning}}$   
2 Right tread speed:  $ts_{\text{right}} = s_t \cdot s_{\max}^{\text{forward}} - w_t \cdot s_{\max}^{\text{turning}}$   
**Output:** Left tread speed  $ts_{\text{left}}$   
Right tread speed:  $ts_{\text{right}}$

---

## Reward Function

## Neural Networks

## Data Analysis

### 4.1.3 Reinforcement learning algorithms implementations

# Chapter 5

## Experimental results (WIP)

TODO (Macaluso P.):

- Introduction: arguments of the chapter and overview of final results
- Experimental Methodology
  - Preliminaries: experiments with Pendulum-v0
  - Real Life experiments with Cozmo
  - Hyper-Parameters discussion and motivation
  - Algorithms applied and modifications with pseudo-code
- Experiments with Pendulum-v0
  - Comparative analysis between results obtained with DDPG and SAC
  - Results
- Experiments with Cozmo
  - Comparative analysis between results obtained with DDPG and SAC
  - Results



# Chapter 6

## Conclusions (WIP)

TODO (Macaluso P.):

- Comments on the results obtained
- Strong points and weaknesses
- Autocriticism about weaknesses that affect the final result
- Future Work
  - Data efficiency and Model-based approaches through a better study of bibliography
  - Usage of Variational Auto Encoder (VAE): this is a possible additional step towards the creation of the model (VAE used for data generation as a generative model). It must not overwrite the first future work.
  - New Version of Cozmo (Vector) with a better camera and LIDAR (Data fusion);

# Appendix A

## Reinforcement Learning

### A.1 Bellman Equation

The value function is decomposable in the immediate reward  $r_t$  and the discounted state value of the next state. It is possible to obtain the result in eq. (A.1) by writing expectations explicitly.

$$\begin{aligned}
 V^\pi(s) &= \mathbb{E}[g_t | s_t = s] \\
 &= \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s] \\
 &= \mathbb{E}[r_{t+1} + \gamma g_{t+1} | s_t = s] \\
 &= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma \mathbb{E}[g_{t+1} | s_{t+1} = s']] \\
 &= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma V^\pi(s')]
 \end{aligned} \tag{A.1}$$

This equation expresses the relationship between the value of a state and the values of its successor states. It is further possible to derive the Bellman Equation for Action-Value function using the same procedure described above.

The resulting formulas are shown in eq. (2.7) on page 6.

Furthermore, it is possible to obtain the Bellman Equation solution in eq. (A.2) working with matrix notation.

$$\begin{aligned}
 V^\pi &= \mathcal{R}^\pi + \gamma \mathcal{P}^\pi V^\pi \\
 (I - \gamma \mathcal{P}^\pi) V^\pi &= \mathcal{R}^\pi \\
 V^\pi &= (I - \gamma \mathcal{P}^\pi)^{-1} \mathcal{R}^\pi
 \end{aligned} \tag{A.2}$$

---

**Algorithm A.1:** Policy Iteration for estimating  $\pi \sim \pi^*$ 


---

**Input:**  $\pi$  the policy to be evaluated; a small threshold  $\theta$  which defines the accuracy of the estimation

- 1 Initialise  $V(s) \forall s \in \mathcal{S}$  arbitrarily, except that  $V(\text{terminal}) = 0$
- 2  $is\_policy\_stable \leftarrow true$
- 3 **repeat**

- 4 **repeat**

- 5  $\Delta \leftarrow 0$
- 6 **for each**  $s \in \mathcal{S}$  **do**

- 7  $v \leftarrow V(s)$
- 8  $V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r|s, a) [r + \gamma V(s')]$
- 9  $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- 10 **end**
- 11 **until**  $\Delta < \theta$
- 12  $V_\pi \leftarrow V(s)$
- 13 **while**  $true$  **do**

- 14 **for each**  $s \in \mathcal{S}$  **do**

- 15  $old\_action \leftarrow \pi(s)$
- 16  $\pi(s) \leftarrow \arg \max_a \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r|s, a) [r + \gamma V_\pi(s')]$
- 17 **if**  $old\_action \neq \pi(s)$  **then**

- 18  $is\_policy\_stable \leftarrow false$
- 19 **end**
- 20 **end**
- 21 **end**
- 22 **until**  $\neg is\_policy\_stable$

**Output:**  $V^*$  and  $\pi^*$

---

## A.2 Dynamic programming

### Policy iteration algorithm

### Policy improvement theorem

Let  $\pi$  and  $\pi'$  be any pair of deterministic policy such that

$$Q_\pi(s, \pi'(s)) \geq V_\pi(s) \quad \forall s \in S \quad (\text{A.3})$$

Then the policy  $\pi'$  leads to

$$V_{\pi'}(s) \geq V_\pi(s) \quad (\text{A.4})$$

Therefore, the presence of strict inequality in eq. (A.3) for a state leads to a strict inequality of eq. (A.4).

The proof of this theorem is shown in eq. (A.5).

$$\begin{aligned}
 V_\pi(s) &\leq Q_\pi(s, \pi'(s)) \\
 &= \mathbb{E}[r_{t+1} + \gamma V_\pi(s_{t+1}) | s_t = s, a_t = \pi'(s)] \\
 &= \mathbb{E}_{\pi'}[r_{t+1} + \gamma V_\pi(s_{t+1}) | s_t = s] \\
 &\leq \mathbb{E}_{\pi'}[r_{t+1} + \gamma Q_\pi(s_{t+1}, \pi'(s_{t+1})) | s_t = s] \quad (\text{by A.3}) \\
 &= \mathbb{E}_{\pi'}[r_{t+1} + \gamma \mathbb{E}_{\pi'}[r_{t+2} + \gamma V_\pi(s_{t+2}) | s_{t+1}, a_{t+1} = \pi'(s_{t+1})] | s_t = s] \\
 &= \mathbb{E}_{\pi'}[r_{t+1} + \gamma r_{t+2} + \gamma^2 V_\pi(s_{t+2}) | s_t = s] \\
 &\leq \mathbb{E}_{\pi'}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V_\pi(s_{t+3}) | s_t = s] \\
 &\vdots \\
 &\leq \mathbb{E}_{\pi'}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots | s_t = s] \\
 &= v_{\pi'}(s)
 \end{aligned} \quad (\text{A.5})$$

### Value iteration algorithm

---

**Algorithm A.2:** Value Iteration, for estimating  $\pi \sim \pi^*$

---

**Input:** A small threshold  $\theta$  which defines the accuracy of the estimation

- 1 Initialise  $V(s) \forall s \in \mathcal{S}$  arbitrarily, except that  $V(\text{terminal}) = 0$
- 2 **repeat**
- 3      $\Delta \leftarrow 0$
- 4     **for** each  $s \in \mathcal{S}$  **do**
- 5          $v \leftarrow V(s)$
- 6          $V(s) \leftarrow \max_a \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma V(s')]$
- 7          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- 8     **end**
- 9 **until**  $\Delta < \theta$
- 10 Output a deterministic policy,  $\pi \sim \pi^*$ , such that

$$\pi(s) = \arg \max_a \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} P(s', r | s, a) [r + \gamma V(s')]$$

**Output:**  $V^*$  and  $\pi^*$

---