



Máster en ingeniería de telecomunicaciones

Proyecto de fin de máster

**Identificación y recogida de objetos con un brazo
robótico utilizando técnicas de reinforcement
learning**

Autor

Pablo Iglesia Fernández-Tresguerres

Dirigido por
Philippe Juhel

Madrid
Enero 2021

ABSTRACT

Abstract content

Thank yous

And other important information

CONTENTS

1. <i>Introduction</i>	1
1.1 Project Motivation	3
2. <i>Description of Technologies</i>	4
3. <i>State of the art</i>	5
3.0.1 Reinforcement Learning	6
3.0.2 Deep Reinforcement Learning	9
3.0.3 Problems of Deep Reinforcement Learning in Real-world . .	16
4. <i>Definition of the Project</i>	18
4.1 Motivation	18
4.2 Objectives	18
4.3 Methodology	19
4.4 Planning and budget	20
5. <i>Developed System</i>	22
5.1 Análisis del sistema	22
5.2 Diseño	22

5.3	Implementación	22
6.	<i>Results Analysis</i>	23
7.	<i>Conclusions and Future Work</i>	24
	<i>Appendix</i>	25
.1	Robot Controller	26
.1.1	main.py	26
.1.2	Robot.py	28
.2	Artificial Intelligence Manager	34
.2.1	main.py	34
.2.2	RAlgorithm.py	36
.2.3	Environment.py	54
.2.4	ImageController.py	57
	<i>Bibliography</i>	59

LIST OF FIGURES

1.1	Pick and Place Task	2
3.1	Fanuc DNN	6
3.2	Q-Matrix	8
3.3	Deep Q Learning	10
4.1	Methodology	20
4.2	Chronograph HW	20
4.3	Planning AI	21
4.4	Planing Robot Controller	21

LIST OF TABLES

1. INTRODUCTION

Robotics and real life are worlds destined to meet. Today everyone has seen robots trying to behave like human beings. Many of them even look similar to a person and try to imitate the way we walk, talk or, ultimately, interact with the environment around us.

Robots, Artificial Intelligence or other concepts such as Machine Learning have crept into our lives in just a few years. In fact, until recently, only a few visionaries like Marvin Minsky or Isaac Asimov used to speak of these concepts, and it was as part of science fiction novels. Nowadays, series like Black Mirror bring this technologies closer to the general public and make us reflect on how the future could be.

However, robots, and artificial intelligence in general, are still far from the vision that is told in the novels. They are not capable of understanding the environment around them, of learning or generalizing as we humans do. Companies and researchers are working on getting better generalization of the algorithms, but the truth is that, so far, Artificial Intelligence is only able to perform specific tasks for which they are programmed.

This project is one of those cases. The goal is to control a UR3 arm robot using Artificial Intelligence in order to pick disordered objects from a box and place them in a point of delivery. This task seems trivial, because we are used to see machines performing pick and place actions in industrial processes, but in fact, these kind of processes are normally just repeating the same action or the same rule over and over again. They are able to perform this tasks because they know apriori where these objects are or how they are placed, but they are not capable of generalizing the workflow.

For instance, in Universal Robot free e-Learning course [1], they expose the following example of an industry pick and place task. In [Figure 1.1](#) we can see how the robot is placing an object in a box located in a conveyor belt. The robot is using an infrared sensor to know that a box has arrived, and this box will always

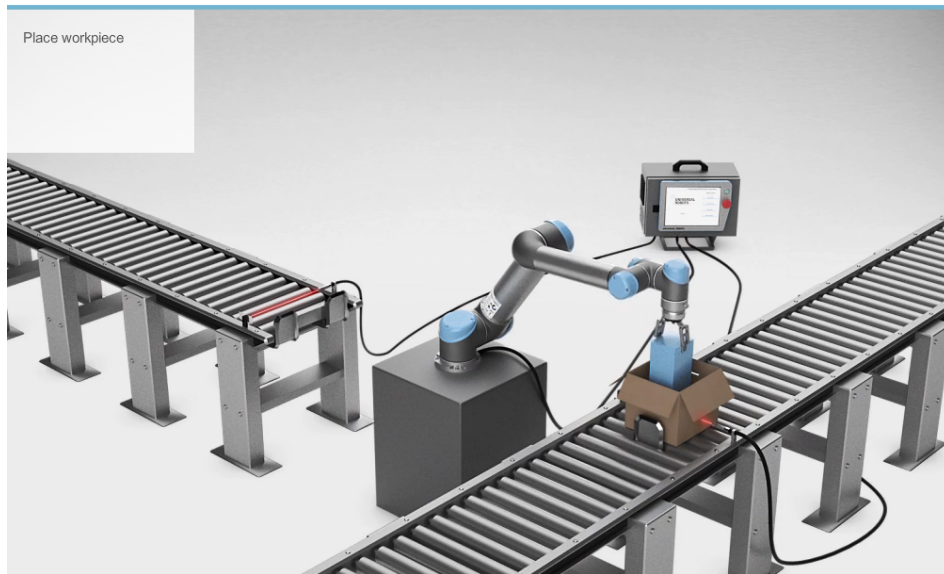


Fig. 1.1: Universal Robots Pick and Place Task

be in the same place because there is a stopper in the conveyor belt which doesn't allow the box to keep moving. On the other hand, the object is picked from the other conveyor belt using the same system to detect the arrival of a new object. The whole task is using a complex architecture, but the robot is performing the same chain of movements in a loop and the only intelligence that the robot has to have is waiting for the object and the box to come.

To achieve generalization in this project, Reinforcement Learning (RL) together with Image Recognition techniques have been used. This algorithms give the robot the ability of calculating, for each time step, the optimal action to achieve the final goal of picking all the objects from the box and placing them in the objective point. To compute this action the robot needs to gather information about the environment such as its relative position over the box or how the pieces are distributed. This information together is called state, and the robot computes each action depending on it.

To perform this project, a distributed architecture with multiple nodes has been created. Each of them takes care of a different activity. For example, some nodes are used to control the robot, others to gather information about the current state, and others are used to train the Artificial Intelligence algorithm. This architecture has been created using ROS (Robot Operative System) and contributes to the

project adding all the advantages of a microservices oriented architecture.

1.1 Project Motivation

The fourth industrial revolution is here, and it will change the way that goods are produced, raising efficiency by increasing the amount of automated processes. This will lead to a faster production and a reduction of errors, as machines have the ability to decide and act in fractions of seconds without making mistakes. Furthermore, machines can also be working 24 hours per day stopping just for maintenance checks, which would help to increase the productivity factor without increasing the expense in human resources.

We have been hearing about industry 4.0 since 2011, but the truth is that it is not a reality yet. We are just in the beginning, and it will take decades to perform such a big change in the industry. There are some factors to take in mind in order to analyse the evolution of the industry in the following years. The improvement on the telecommunications with the arrival of 5G networks, the moral dilemma of substituting workers for machines and the impact that this could have in the society or the improvement and implementation of AI technologies are just some of these factors.

We have seen a lot of Artificial Intelligence algorithms applied to the industry, but the truth is that these technologies are not fully developed yet and just big companies can afford to use them in their supply chain. Besides, there are some task that are now performed by humans and cannot be done by machines due to its complexity or its importance in the whole production chain.

The motivation of this project is to contribute to the industry change providing an open source solution to a complex problem such as disordered pick and place task. This open source solution does not currently exist in the industry and would add value being a good starting point for bigger projects in the future.

2. DESCRIPTION OF TECHNOLOGIES

Describir las tecnologías, protocolos, herramientas específicas, etc. que se vayan a tratar durante el proyecto para facilitar su lectura y comprensión. Hablar de Java no procede aquí porque todo el mundo sabe lo que es, pero si en el proyecto hablo continuamente del protocolo Baseband, debo especificar en este capítulo qué es y para qué sirve.

- ROS + catkin
- pytorch
- arduino
- github
- CUDA
- moveit
- UR driver
- UR3 robot
- gazebo
- arduino
- Reinforcement learning
- anaconda

3. STATE OF THE ART

The pick and place task that is intended to be performed in this thesis is really useful for a lot of applications into the industrial world because it would bring a lot of flexibility for these processes. A example of this applications could be an assembly line, where robotic arms could be picking all the different pieces to assemble in the product using always the same algorithm.

Big companies are developing a lot of Artificial intelligence use cases in the industry, and they try to contribute to the AI community by publishing scientific articles on how they managed to use AI for their specific tasks. Unfortunately, although some companies have already developed their own solutions for our specific pick and place task, none of them have published a scientific article on the subject, making it difficult to study the way they have achieved it.

One of the companies that has already developed a pick and place task is the Japanese automation company Fanuc, which has developed an AI-based solution together with Preferred Networks. As commented before, they have not published any scientific article about the topic but we can see the system working in a video they have posted on YouTube [2]. That means that we have to gather all the possible information from the video, where we could find that they have not used a Reinforcement Learning algorithm but just a Deep Neural Network (DNN) with image recognition.

To train the net, they have collected "success" or "fail" labelled images by making pick actions in random places of the box. Once they gathered a big enough dataset of images, they have trained a Deep Neural Network as the one we see in [Figure 3.1](#), where we can also see that the Neural Net has been trained to predict whether the robot is going to success in a pick action in a specific place or not. Using that net, they can make a heat map of the whole box, predicting the points of maximum probabilities of succeed. As usually, they noticed that the bigger the image dataset, the higher the success ratio. In eight hours, they reach 90% of success, which they say is bigger than the human success ratio.

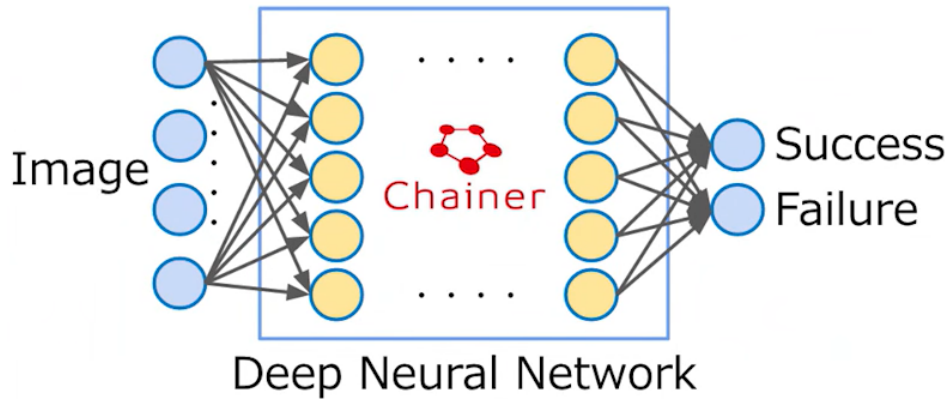


Fig. 3.1: Deep Neural Network of Fanuc solution (taken from the video)

3.0.1 Reinforcement Learning

The idea of the project is to keep using image recognition techniques but, in our case, applied to a **Reinforcement Learning** Algorithm which is an area of machine learning inspired by psychology behavioural. Its goal is to determine what actions a software agent should choose in a given environment in order to maximize some notion of "reward" or accumulated prize.

Explained easily, RL is used to make an **agent** (the robot) learn how to interact with a **environment** in order to perform a task. To achieve this, Markov Decision Process (MDP) which provides a mathematical framework for modeling decision making in situations where outcomes are partly random and partly under the control of a decision maker.

Markov Decision Process (MD)

In MDP, the environment is what we are actually trying to simulate with the MDP. The agent will interact with it to learn how to perform the task, so these are the attributes of the environment:

- **Agent:** The agent is the most important piece of the algorithm because it represents the objects that we want to become smarter.
- **Actions (A):** The agent can interact with the environment by performing a set of actions which is normally finite.

-
- **States (S):** Each time the agent performs an action, it moves to a new state. States are basically the set of information that differentiates the situation of the agent before and after performing an action. States can be transitional or terminal, when the agent meets the objective or when it gets to a forbidden position.
 - **Rewards (R):** Each time an action is performed, the agent receives a reward. This reward can be positive, negative or null depending on the impact of the action to achieve the objective.
 - **Policy (π):** The policy is used to define the optimal action for each step. It gives a punctuation for all of the actions in the current step as shown in the following formula. The agent takes the action with highest punctuation.

$$\pi(a|s) = P_r\{A_t = a|S_t = s\}$$

The MDP is divided in discrete timesteps (t), where each timestep does not have to last the same time as the previous step. Each timestep, the agent uses the policy π to decide the next action.

Once the action is taken:

- The environment transits to th next state: $S_t = S_{t+1}$.
- Environment produces a new reward, which can be represented with the following formula:

$$P(s', r, s, a) = P_r\{S_{t+1} = s', R_t = r, S_t = s, A_t = a\}$$

Agent's performance is calculated in terms of its future accumulated rewards known as return. This is called **expected return** an is calculated as shown in the formula below, where γ is the discount factor, and is used to give a bigger value to the closest steps.

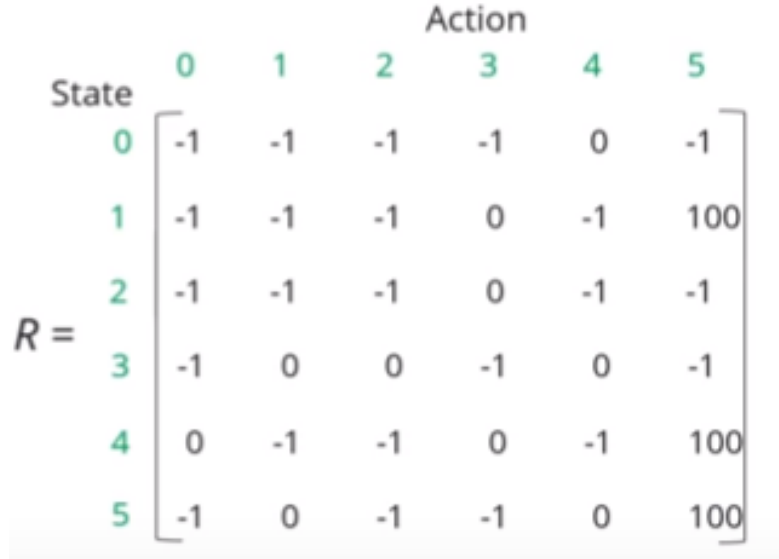
$$G_t = \sum_{k=t} \gamma^{k-t} \cdot R_{k+1} \quad \forall \gamma \in [0, 1]$$

Q-Learning

Now that we know all these concepts, we have to learn what Reinforcement Learning Algorithms do to learn. Basically, **the goal of the agent is to find a policy that maximizes the expected return**. This can be done using different strategies as:

- **Q-Learning:** Estimating action values using Q Tables or other methods
- **TRPO:** Parametrizing the policy and optimizing its parameters

Basic Q-Learning is based on the assumption that both actions and states are limited and that the same action in the same state always drives to the same new state. Having this in mind, Q-learning algorithms build two matrices of shape $\text{length}(\text{actions}) \times \text{length}(\text{states})$ as shown in the [Figure 3.2](#).



		Action					
		0	1	2	3	4	5
State	0	-1	-1	-1	-1	0	-1
	1	-1	-1	-1	0	-1	100
	2	-1	-1	-1	0	-1	-1
	3	-1	0	0	-1	0	-1
	4	0	-1	-1	0	-1	100
	5	-1	0	-1	-1	0	100

Fig. 3.2: Reward and Q Matrix shape in Basic Reinforcement Learning

In these two matrices, Q-Learning algorithm stores in the R matrix the reward for the pair of action-state while in the Q matrix they store cumulated reward for this same pair. The Q matrix is the one used to decide which action to perform in each state and R matrix the one used to calculate the reward of each action.

However, for the aim of this project, the states of the agent can be different in each timestep. The state would actually be partially formed by images, so the

number of states can be infinite. We need a more complex version of Reinforcement Learning.

3.0.2 Deep Reinforcement Learning

The approach of mixing both image recognition and RL is called Deep Q Learning (DQN) or Double Deep Q Learning (DDQN) depending on the implementation and uses Neural Networks in two different stages of the algorithm. Firstly, a Convolutional Neural Network (CNN) is used to extract image features, and then, a Deep Neural Network (DNN) is used to calculate the q value of each independent action and select the next one using these values.

DQL was proposed in 2012, and, since then, it has been used for a lot of different purposes. For example, Guillaume Lample and Devendra Singh Chaplot demonstrated back in 2017 that a RL agent could play FPS Games using as inputs just game scores and pixels from the screen [3]. Another really interesting example is this robot [4], which is capable of moving around a house looking for an objective and avoiding obstacles using DDQL.

A good resource to understand how Reinforcement Learning really works is Deeplizard's tutorial [5]. In this tutorial they explain different versions of the algorithm and how to implement them in python to solve different OpenAI gym environments [6].

Deep Reinforcement Learning is though an union between RL and image recognition, but let's see how it actually works. The main idea is to replace the Q-table that we saw before for a Dense Neural Network that uses as input another Neural Network, a Convolutional Neural Network (CNN). The full algorithm would have as many outputs as allowed actions. Therefore, simplifying, these outputs are equivalent to the q-values saw before and so we will call them. To see it graphically, when the agent wanted to take an action, he would pass the state image through the Neural network represented in [Figure 3.3](#) and would take the action with higher q-value.

When I said "simplifying" in the previous paragraph, I meant "simplifying a lot" in the next paragraphs I will explain all the intermediate steps in the algorithm and why they are important:

- Episodes and Steps

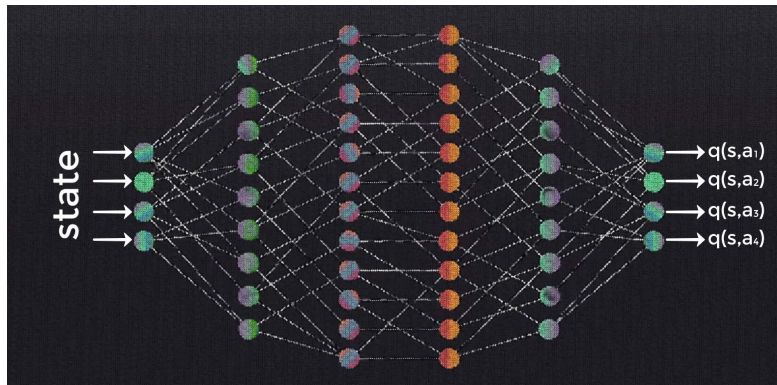


Fig. 3.3: Deep Q Learning Representation with 4 outputs

- Exploration vs Exploitation trade-off
- Replay Memory
- Bellman's Equation
- Target and Policy Networks

Episodes and Steps

RL training is divided in Episodes. One Episode is the sequence of actions needed to reach a terminal State. Each time the agent reaches a terminal state, an episode is ended, and a new one is started.

On the other hand, steps represents every time that a new action is taken, so the number of steps taken by the agent during training is infinite. Later on, we will use as metric of performance the number of steps per episode, as they must decrease during the training.

Exploration vs Exploitation trade-off

In Reinforcement Learning there are two important concepts that are **Explore** and **Exploit**. To explore is basically gather new information about the environment and to exploit is to make the best decision with the information that we already have.

In Deep Reinforcement Learning, the agent exploit the information gathered by using the pre-trained Neural Network to decide next action. On the other hand, the agent explore the environment by deciding next action randomly. We use exploration mainly in the beginning of the training because we want the agent to gather as much information of the environment as possible before starting training.

When the agent uses exploitation, it is also gathering information about the environment. However, we could not let the agent explore this way because during the exploration phase we want all the actions to be performed with the same probability and neural network bias can cause some actions to be performed much more than others.

So, how do we decide when the agent must explore or exploit? To decide it we can use multiple techniques, but the most common one is the Epsilon-Greedy Strategy. This strategy basically consist on setting a probability of exploring and keep decreasing it slowly during the training. It works this way:

1. We set the initial exploring probability (ϵ)
2. We set the per-step epsilon decay, (ϵ_{decay})
3. For each step:
 - (a) With probability $p = \epsilon$, the agent explores the environment (takes a random action). If not, it exploit the information by deciding the action using the NN.
 - (b) Whether the agent has explore or not, we decrease the probability of exploring the environment in the next step ($\epsilon = \epsilon - \epsilon_{decay}$)

Using this strategy, the agent will rather explore or exploit the environment during the training. In the first steps the probability of a random action (exploring) will be much hihger than in the last steps of the algorithm. This probability will keep decreasing during the training, until it reaches the minimum exploring rate, which is normally set to 10

Replay Memory

Every time that the agent performs an action, either by exploring or exploiting, the agent lives an experience. For the purpose of training the algorithm, we will store all these experiences.

Experiences are formed by the initial state, the action taken, the state reached (final state) and the reward gotten and they are stored in the Replay Memory. Then, every time that an action is taken, the algorithm is trained following this steps:

1. Replay Memory checks if the number of experiences is higher than the batch size
2. If there are enough experiences:
 - (a) Replay Memory supplies a random set of experiences of size=batch_size.
 - (b) With this set of experiences, the target network is trained.

Optimizing Replay Memory can be a challenge, because, if we are using a Graphic Card in the training, we would be storing all the experiences in its memory. But, why do we need to store all the experiences? We could also be using the last N experiences to train the network and it would be a less memory-consumption demanding solution.

The answer to this question is that Reinforcement Learning Networks converge really slowly and variance between consecutive steps is really low. Using consecutive experiences to train the network would result though in a slower and biased learning. Besides, this way of working is better for learning real-world experience, where there are infinite different states, as the experience gained in previous steps will be used multiple times later to train the network.

Bellman's Equation

As commented before, Deep Reinforcement Learning uses Neural Networks to compute the q-values of each action. The optimal value of these q-values is represented by the Bellman's Equation and is shown below:

$$q_*(s, a) = E[R_t + \gamma \max(q_*(s', a'))]$$

As we can see in the equation, the optimal value depends in both the reward of the action taken and the maximum optimal q-value of the next action. In real life it is impossible to compute this value, because we would be an infinite loop. However, as the most important parameter of the formula is the expected reward ($q_*(s', a')$ is multiplied by the discount factor γ), we can simply use the next action q-value and it will be a good approximation. The formula would stay as follows:

$$q_*(s, a) = E[R_t + \gamma \max(q(s', a'))]$$

With this new formula we will be able to compute the optimal q-value for each experience stored in the Replay Memory (initial state, action, final state and reward). It is important to have in mind for this process that the optimal q-value can only be computed if the action has actually been taken, because we don't know the Reward of non taken actions. But, anyway, why do we need to compute the optimal q-value?

To answer this question, lets take a look to the training process of the neural network:

1. the agent decide which action to take using the policy-network. (action with highest q-value)
2. The agent takes the action and receives a reward from the environment
3. The agent stores all the experience in the Replay Memory
4. The training process is started:
 - (a) A random batch of experiences is taken from the Replay Memory
 - (b) For all these experiences, its optimal q-value is calculates using the modified Bellman's Equation and target-network
 - (c) For all these experiences, the actual q-value is calculated using the policy-network
 - (d) For all these experiences, the loss is calculated as the difference of both values

-
- (e) We use the Neural network optimizer to back-propagate the loss to all the weights

So, to answer the previous question, we need to compute the **optimal q-values in order to calculate the loss** of the neural network for each action taken and train, though, the algorithm.

Retaking here the question answered before about why we needed Replay Memory module, one important reason is that one action taken in the initial steps of the training will affect differently to the neural network in this moment than later, when the network is already trained, and its q-value is though more similar to the optimal q-value. Replay Memory technique allow us to use this information gathered in any step of the training, during a step where the network is more trained.

Target and Policy Networks

In the previous step, we talk about two different networks: policy and target. The target network comes to solve a stability problem of the DRL training. In the next paragraphs I will explain the problem and how target network can help to solve it.

Having in mind the way we calculate the loss of the neural network in the previous section we can realize that we have to pass information through the network twice. Just to remember:

$$loss = R_t + \gamma max(q(s_{t+1}, a_{t+1})) - q(s_t, a_t)$$

As a spoiler, I can say that $q(s_{t+1})$ and $q(s_t, a_t)$ will not be calculated with the same network. But.. why?

Let's see what could happen if we calculated both of the values with the same network. If we used the same network, we would calculate the loss of the neural network as the difference between two consecutive

We do the first pass to calculate the Q-value for the relevant action, and then we do a second pass in order to calculate the target Q-value for this same action. Our objective is to get the Q-value to approximate the target Q-value.

Remember, we don't know ahead of time what the target Q-value even is, so we attempt to approximate it with the network. This second pass occurs using the same weights in the network as the first pass.

Given this, when our weights update, our outputted Q-values will update, but so will our target Q-values since the targets are calculated using the same weights. So, our Q-values will be updated with each iteration to move closer to the target Q-values, but the target Q-values will also be moving in the same direction.

As Andong put it in the comments of the last video, this makes the optimization appear to be chasing its own tail, which introduces instability. As our Q-values move closer and closer to their targets, the targets continue to move further and further because we're using the same network to calculate both of these values.

Well, here's a perfect time to introduce the second network that we mentioned earlier. Rather than doing a second pass to the policy network to calculate the target Q-values, we instead obtain the target Q-values from a completely separate network, appropriately called the target network.

The target network is a clone of the policy network. Its weights are frozen with the original policy network's weights, and we update the weights in the target network to the policy network's new weights every certain amount of time steps. This certain amount of time steps can be looked at as yet another hyperparameter that we'll have to test out to see what works best for us.

So now, the first pass still occurs with the policy network. The second pass, however, for the following state occurs with the target network. With this target network, we're able to obtain the Q-value for the next state, and again, plug this value into the Bellman equation in order to calculate the target Q-value for the first state.

This is all we use the target network for — To find the value of this term so that we can calculate the target Q-value for any state passed to the policy network. As it turns out, this removes much of the instability introduced by using only one network to calculate both the Q-values, as well as the target Q-values. We now have something fixed, i.e. fixed Q-targets, that we want our policy network to approximate. So, we no longer have the dog-chasing-it's-tail problem.

As mentioned though, these values don't stay completely fixed the entire time. After amount of time steps, we'll update the weights in the target network with the weights from our policy network, which will in turn, update the target Q-values with respect to what it's learned over those past time steps. This will cause the

policy network to start to approximate the updated targets.

3.0.3 Problems of Deep Reinforcement Learning in Real-world

Real-world problems introduces some challenges that we will have to manage. In march 2018, A. Rupam Mahmood, Dmytro Korenkevych, Brent J. Komer, and James Bergstra explained the problems they found while implementing a RL algorithm in a UR5 robotic arm [7].

Some of the problems they found were the following:

- Slow rate of data-collection, as movements in the real robot are slower than in a simulated environment.
- Partial observability. Sensors cannot retrieve all the information about the environment.
- Noisy sensors will provide inaccurate information.
- Safety of the robot and its surroundings have to be taken in mind.
- Fragility of robot components.
- Delay between an action is requested and the time it is actually performed can affect the training.
- Preparing the robot is a really difficult task:
 - Controlling the robot.
 - Define all aspects of the environment.
 - Difficulties for obtaining random and independent state when episode ends.

Another problem that can be found in our project is that, as objects are randomly placed, the environment that the agent will have to face will be completely different each time. In fact, the robot can interact with the environment, as it can move the pieces trying to pick them, so we are facing a dynamic environment RL problem. A good example of a dynamic environment problem is the path planning of a self driven car, where each time the agent takes an action the environment

will change and, furthermore, obstacles do not have to be static, but they can also move.

There are multiple examples of articles on this topic, such as the one Xiaoyun Lei, Zhian Zhang, and Peifang Dong published in September 2018 using a DDQN approach to solve it [8]. However, there are other solutions as the one proposed by Marco A. Wiering [9], where he introduces some prior knowledge to the model in order to facilitate the learning. His algorithm had problems generalizing the environment, so he introduced some prior information about the model together with a Model-based RL. This made the algorithm more capable to learn without losing a lot of trainable capability.

4. DEFINITION OF THE PROJECT

4.1 *Motivation*

The project motivation is the natural continuation of a previous project performed at ICAM University. This project was part of the assembly line of a car manufacturing process and its objective was to pick some specific plastic pieces and place them into the product. To achieve this, the system used opencv image processing, so it was recognising a specific shape given apriori.

This project was totally functional and the robot could perform the task with a high successful rate. However, the lack of generalization of the system makes it hard to introduce changes as using it for another part of the assembly line. Each time that this happened someone would have to introduce the shape of the pieces to the system and to calibrate the camera to the new environment.

The motivation of the project is to create from scratch a new solution for performing the picking of the pieces. This time, the project will not be sponsored by any company, so there will be less resources to use.

With this new approach, the idea is to use all the knowledge of previous documented projects on Artificial Intelligence in the industry and make a little contribution to the huge advance of industry 4.0. In fact, the idea is to make the project completely replicable so that anyone could use this project as a starting point for new applications.

4.2 *Objectives*

The objectives of the project are five and are listed and explained below:

- Implement a bin picking simple solution. A basic one, without Artificial

Revisar
obje-
tivos
y
adap-
tar-
los al
proyecto

Intelligence.

- Improve the performance using RL and Image Recognition.
- Study the usage of new technologies to add information to the system. Adding physics or new image recognition techniques.
- Test different tools such as Multiflash or 3D images to improve performance.
- Add value to the scientific community making the solution available if the previous objectives are reached.

4.3 Methodology

This project will be performed using an agile methodology, which is one of the simplest and effective processes to turn a vision for a business need into software solutions. Agile is a term used to describe software development approaches that employ continual planning, learning, improvement, team collaboration, evolutionary development, and early delivery. It encourages flexible responses to change [10].

In the case of this project, the team is just formed by two workers and a project manager. This make necessary to make some changes to the typical agile methodology. For example, daily meetings are substituted by constant communication between both workers and weekly meetings with the project director. With this approach, all the members of the project are updated about how it is going and have clear objectives.

Likewise, the methodology of this project is based in three fundamental principles as it is shown in [Figure 4.1](#). The first two principles are highly related, as the project is iterative because it is experimental. That means that the way of working is perform little sprints with new functionalities, test them (experimental) and, depending on the results, define the new sprint, execute it and test again (iterative). Besides, the project is also incremental because the idea is starting implementing the simplest possible solution and keep adding new improvements to it in a iterative loop until the optimal configuration is reached.

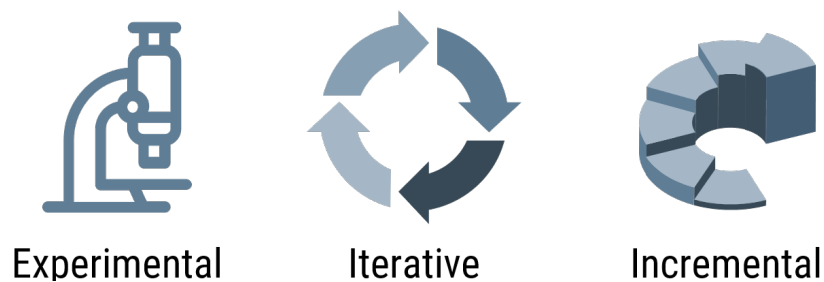


Fig. 4.1: Methodology

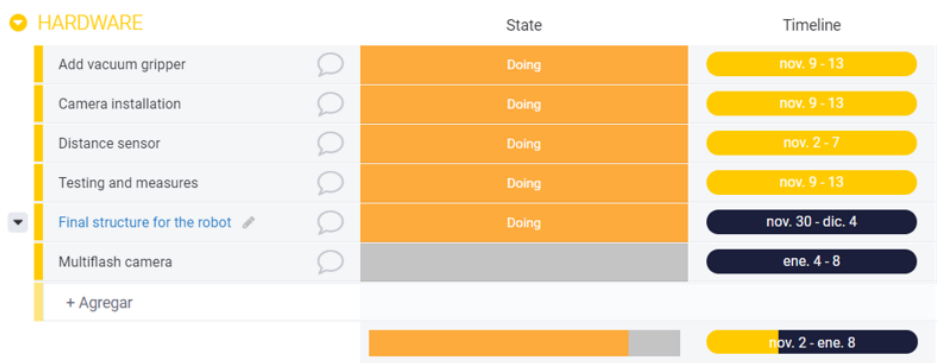


Fig. 4.2: Chronograph of the Hardware implementation

4.4 Planning and budget

Regarding the planning, there are some really important functionalities that have been defined since the beginning of the project, as they are needed. These functionalities are split in three different groups: Hardware implementation, Artificial Intelligence Implementation and Robot controller implementation. The tasks related to these three groups are shown in [Figure 4.2](#), [Figure 4.3](#) and [Figure 4.4](#).

Añadir pre-supuesto

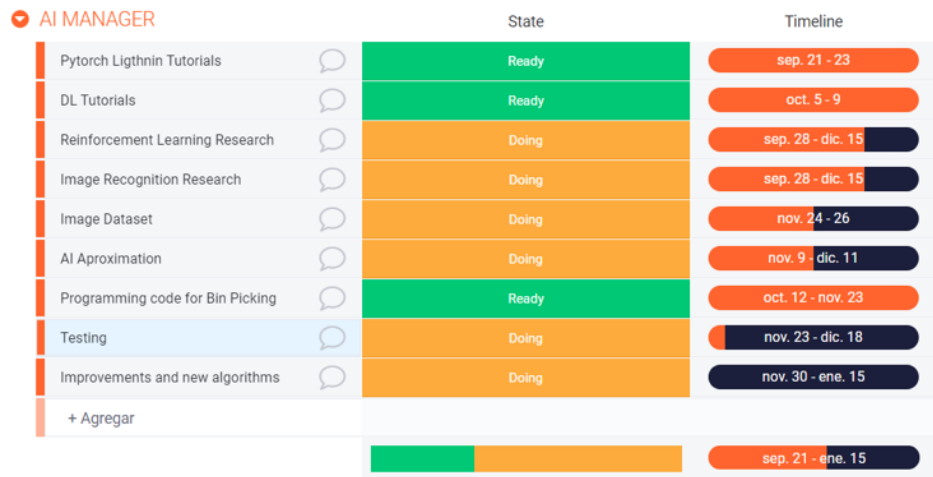


Fig. 4.3: Planning of the Artificial Intelligence implementation

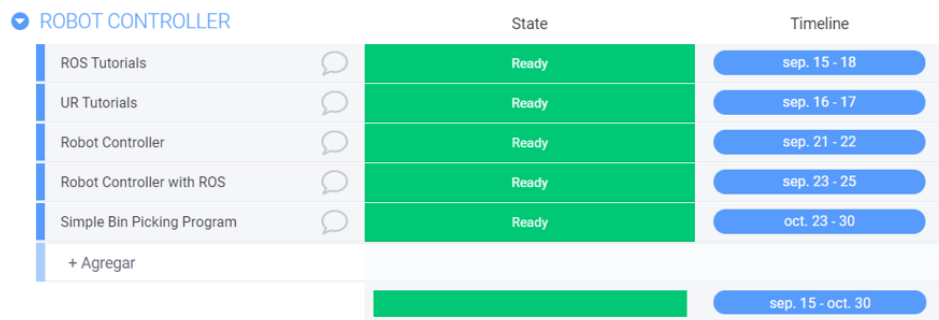


Fig. 4.4: Planning of the Robot Controller implementation

5. DEVELOPED SYSTEM

En este capítulo es donde el alumno debe describir su proyecto. En función del tipo de proyecto la estructura interna variará. El título del mismo, así como sus apartados, son sólo una sugerencia que cada alumno deberá adaptar particularmente a su proyecto.

A su vez, este capítulo podrá extenderse en varios capítulos más, por ejemplo si mi proyecto consta de varias fases o módulos, lo lógico sería tener varios capítulos donde se describa el desarrollo del proyecto:

- Capítulo 5. Implantación y configuración de la plataforma
- Capítulo 6. Desarrollo del sistema

5.1 Análisis del sistema

...

5.2 Diseño

...

5.3 Implementación

...

6. RESULTS ANALYSIS

Destacar los resultados más relevantes del proyecto y hacer un análisis crítico de los mismos. También es un capítulo obligatorio y clave.

7. CONCLUSIONS AND FUTURE WORK

Comentar las conclusiones del proyecto, destacando lo que se ha hecho, dejando claros qué objetivos se han cubierto y cuáles son las aportaciones hechas.

APPENDIX

.1 Robot Controller

En esta sección se muestra el código de algunos de los principales elementos del módulo Robot Controller, implementado en la arquitectura del proyecto.

Este código también está disponible en el siguiente repositorio de github:

https://github.com/pabloiglesia/robot_controller

.1.1 main.py

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 """
5 - We need to connect the camera and the nodes
6 roslaunch ur_icam_description webcam.launch
7
8 - We need to establish a connection to the robot with the
  following comand:
9 roslaunch ur_robot_driver ur3_bringup.launch robot_ip
  :=10.31.56.102 kinematics_config:=${HOME}/Calibration/
  ur3_calibration.yaml
10
11 - Then, we ned to activate moovit server:
12 roslaunch ur3_moveit_config ur3_moveit_planning_execution.launch
13
14 - Activate the talker
15 rosrun ai_manager main.py
16
17 - Activate the node
18 rosrun robot_controller arduino.py
19
20 - Finally, we can run the program
21 rosrun robot_controller main.py
22
23 """
24
25 import rospy
26
27 from ai_manager.srv import GetActions
28 from Robot import Robot
29
30
31 def get_action(robot, object_gripped):
```

```
32     relative_coordinates = robot.calculate_current_coordinates()
33     rospy.wait_for_service('get_actions')
34     try:
35         get_actions = rospy.ServiceProxy('get_actions', GetActions
36     )
37         return get_actions(relative_coordinates[0],
38         relative_coordinates[1], object_gripped).action
39     except rospy.ServiceException as e:
40         print("Service call failed: %s"%e)
41
42 # This function defines the movements that robot should make
43 # depending on the action listened
44 def take_action(action, robot):
45     rospy.loginfo("Action received: {}".format(action))
46     object_gripped = False
47     if action == 'north':
48         robot.take_north()
49     elif action == 'south':
50         robot.take_south()
51     elif action == 'east':
52         robot.take_east()
53     elif action == 'west':
54         robot.take_west()
55     elif action == 'pick':
56         object_gripped = robot.take_pick()
57     elif action == 'random_state':
58         robot.take_random_state()
59     return object_gripped
60
61 if __name__ == '__main__':
62     rospy.init_node('robotUR')
63
64     robot = Robot()
65
66     # Test of positioning with angular coordinates
67     robot.go_to_initial_pose()
68
69     # Let's put the robot in a random position to start, creation
70     # of new state
71     object_gripped = take_action('random_state', robot)
72
73     while True:
74         action = get_action(robot, object_gripped)
75         object_gripped = take_action(action, robot)
```

.1.2 Robot.py

```
1 import copy
2 import rospy
3 import time
4
5 from std_msgs.msg import Bool
6 from std_msgs.msg import Float32
7
8 from ai_manager.Environment import Environment
9 from ur_icam_description.robotUR import RobotUR
10
11 """
12 Class used to establish connection with the robot and perform
13 different actions such as move in all cardinal directions
14 or pick and place an object.
15 """
16
17 class Robot:
18     def __init__(self, robot=RobotUR(), gripper_topic='
switch_on_off'):
19         self.robot = robot # Robot we want to control
20         self.gripper_topic = gripper_topic # Gripper topic
21         self.gripper_publisher = rospy.Publisher(self.
gripper_topic, Bool) # Publisher for the gripper topic
22
23     def relative_move(self, x, y, z):
24         """
25         Perform a relative move in all x, y or z coordinates.
26
27         :param x:
28         :param y:
29         :param z:
30         :return:
31         """
32         waypoints = []
33         wpose = self.robot.get_current_pose().pose
34         if x:
35             wpose.position.x -= x # First move up (x)
36             waypoints.append(copy.deepcopy(wpose))
37         if y:
38             wpose.position.y -= y # Second move forward/backwards
39         in (y)
40             waypoints.append(copy.deepcopy(wpose))
41         if z:
42             wpose.position.z += z # Third move sideways (z)
43             waypoints.append(copy.deepcopy(wpose))
```

```
44         self.robot.exec_cartesian_path(waypoints)
45
46     def calculate_relative_movement(self, relative_coordinates):
47         absolute_coordinates_x = Environment.CARTESIAN_CENTER[0] -
relative_coordinates[0]
48         absolute_coordinates_y = Environment.CARTESIAN_CENTER[1] -
relative_coordinates[1]
49
50         current_pose = self.robot.get_current_pose()
51
52         x_movement = current_pose.pose.position.x -
absolute_coordinates_x
53         y_movement = current_pose.pose.position.y -
absolute_coordinates_y
54
55         return x_movement, y_movement
56
57     def calculate_current_coordinates(self):
58         absolut_coordinate_x = self.robot.get_current_pose().pose.
position.x
59         absolut_coordinate_y = self.robot.get_current_pose().pose.
position.y
60
61         relative_coordinate_x = Environment.CARTESIAN_CENTER[0] -
absolut_coordinate_x
62         relative_coordinate_y = Environment.CARTESIAN_CENTER[1] -
absolut_coordinate_y
63
64         return [relative_coordinate_x, relative_coordinate_y]
65
66     # Action north: positive x
67     def take_north(self, distance=Environment.ACTION_DISTANCE):
68         self.relative_move(distance, 0, 0)
69
70     # Action south: negative x
71     def take_south(self, distance=Environment.ACTION_DISTANCE):
72         self.relative_move(-distance, 0, 0)
73
74     # Action east: negative y
75     def take_east(self, distance=Environment.ACTION_DISTANCE):
76         self.relative_move(0, -distance, 0)
77
78     # Action west: positive y
79     def take_west(self, distance=Environment.ACTION_DISTANCE):
80         self.relative_move(0, distance, 0)
81
82     def take_random_state(self):
83         # Move robot to random positions using relative moves. Get
coordinates
```

```
84         relative_coordinates = Environment.generate_random_state()
85         # Calculate the new coordinates
86         x_movement, y_movement = self.calculate_relative_movement(
relative_coordinates)
87         # Move the robot to the random state
88         self.relative_move(x_movement, y_movement, 0)
89
90     def send_gripper_message(self, msg, timer=2, n_msg=10):
91         """
92         Function that sends a burst of n messages of the
gripper_topic during an indicated time
93         :param msg: True or False
94         :param time: time in seconds
95         :param n_msg: number of messages
96         :return:
97         """
98         time_step = (timer/2)/n_msg
99         i=0
100         while(i <= n_msg):
101             self.gripper_publisher.publish(msg)
102             time.sleep(time_step)
103             i += 1
104
105         time.sleep(timer/2)
106
107     # Action pick: Pick and place
108     def take_pick(self):
109         # In this function we should read the distance to the
object
110         # up_distance = 0 # Variable were we store the distance
that we have move the robot so that we can go back to the
111         # original pose
112
113         def change_plan_speed(plan, new_speed):
114             """
115             Function used for changing Robot velocity of a
cartesian path once the movement have been planned.
116             :param plan: RobotTrajectory object. For example, the
one calculated by compute_cartesian_path() MoveGroup function.
117             :param new_speed: speed factor of the robot, been 1
the original speed and 0 the minimum.
118             :return: RobotTrajectory object (new plan).
119             """
120             new_plan = plan
121             n_joints = len(plan.joint_trajectory.joint_names)
122             n_points = len(plan.joint_trajectory.points)
123
124             points = []
125             for i in range(n_points):
```

```
126         plan.joint_trajectory.points[i].time_from_start =
plan.joint_trajectory.points[
127
128             i].time_from_start / new_speed
129             velocities = []
130             accelerations = []
131             positions = []
132             for j in range(n_joints):
133                 velocities.append(plan.joint_trajectory.points
[i].velocities[j] * new_speed)
134                 accelerations.append(plan.joint_trajectory.
points[i].accelerations[j] * new_speed)
135                 positions.append(plan.joint_trajectory.points[
i].positions[j])
136
137             point = plan.joint_trajectory.points[i]
138             point.velocities = velocities
139             point.accelerations = accelerations
140             point.positions = positions
141
142             points.append(point)
143
144         new_plan.joint_trajectory.points = points
145
146         return new_plan
147
148     def back_to_original_pose(robot):
149         """
150         Function used to go back to the original height once a
vertical movement has been performed.
151         :param robot: robot_controller.Robot.py object
152         :return:
153         """
154         distance = Environment.CARTESIAN_CENTER[2] - robot.
robot.get_current_pose().pose.position.z
155         robot.relative_move(0, 0, distance)
156
157     def down_movement(robot, movement_speed):
158         """
159         This function performs the down movement of the pick
action.
160
161         It creates an asynchronous move group trajectory
planning. This way the function is able to receive distance
messages while the robot is moving and stop it once
the robot is in contact with an object.
162
163         Finally, when there is any problems with the
communications the movement is stopped and
```

```
164         communication_problem boolean flag is set to True. It
165         is considered that there is a problem with
166         communications when the robot is not receiving any
167         distance messages during 200 milli-seconds (timeout=0.2)
168
169         :param robot: robot_controller.Robot.py object
170         :return: communication_problem flag
171         """
172
173         distance_ok = rospy.wait_for_message('distance', Bool)
174     .data # We retrieve sensor distance
175         communication_problem = False
176
177         if not distance_ok: # If the robot is already in
178         contact with an object, no movement is performed
179             waypoints = []
180             wpose = robot.robot.get_current_pose().pose
181             wpose.position.z -= (wpose.position.z) # Third
182             move sideways (z)
183             waypoints.append(copy.deepcopy(wpose))
184
185             (plan, fraction) = robot.robot.move_group.
186             compute_cartesian_path(
187                 waypoints, # waypoints to follow
188                 0.01, # eef_step
189                 0.0) # jump_threshold
190
191             plan = change_plan_speed(plan, movement_speed)
192             robot.robot.move_group.execute(plan, wait=False)
193
194             while not distance_ok:
195                 try:
196                     distance_ok = rospy.wait_for_message('
197                     distance', Bool, 0.2).data # We retrieve sensor distance
198                 except:
199                     communication_problem = True
200                     rospy.loginfo("Error in communications,
201                     trying again")
202                     break
203
204             # Both stop and 10 mm up movement to stop the
205             robot
206             robot.robot.move_group.stop()
207             robot.relative_move(0, 0, 0.001)
208
209             return communication_problem
210
211         communication_problem = True
212         while communication_problem: # Infinite loop until the
```

```
movement is completed
204         communication_problem = down_movement(self,
movement_speed=0.2)
205
206         self.send_gripper_message(True, timer=4) # We turn on the
gripper
207
208         back_to_original_pose(self) # Back to the original pose
209
210         object_gripped = rospy.wait_for_message('object_gripped',
Bool).data
211         if object_gripped: # If we have gripped an object we
place it into the desired point
212             self.take_place()
213         else:
214             self.send_gripper_message(False) # We turn off the
gripper
215
216         return object_gripped
217
218     # Function to define the place for placing the grasped objects
219     def take_place(self):
220         # First, we get the cartesian coordinates of one of the
corner
221         x_box, y_box = Environment.get_relative_corner('se')
222         x_move, y_move = self.calculate_relative_movement([x_box,
y_box])
223         # We move the robot to the corner of the box
224         self.relative_move(x_move, y_move, 0)
225         # We calculate the trajectory for our robot to reach the
box
226         trajectory_x = self.robot.get_current_pose().pose.position
.x - Environment.PLACE_CARTESIAN_CENTER[0]
227         trajectory_y = self.robot.get_current_pose().pose.position
.y - Environment.PLACE_CARTESIAN_CENTER[1]
228         trajectory_z = - Environment.CARTESIAN_CENTER[2] +
Environment.PLACE_CARTESIAN_CENTER[2]
229         # We move the robot to the coordinates desired to place
the object
230         self.relative_move(0, 0, trajectory_z)
231         self.relative_move(0, trajectory_y, 0)
232         self.relative_move(trajectory_x, 0, 0)
233         # Then, we left the object
234         self.relative_move(0, 0, -0.05)
235         # Then, we switch off the vacuum gripper so the object can
be placed
236         self.send_gripper_message(False)
237         # Wait some seconds, in order to the msg to arrive to the
gripper
```



```
238         time.sleep(2)
239         # Then the robot goes up
240         self.relative_move(0, 0, 0.05)
241         # Final we put the robot in the center of the box, the
episode should finish now
242         self.robot.go_to_joint_state(Environment.ANGULAR_CENTER)
243
244     def go_to_initial_pose(self):
245         target_reached = self.robot.go_to_joint_state(Environment.
ANGULAR_CENTER)
246
247         if target_reached:
248             print("Target reachead")
249         else:
250             print("Target not reached")
```

.2 Artificial Intelligence Manager

En esta sección se muestra el código de algunos de los principales elementos del módulo AI Manager, implementado en la arquitectura del proyecto.

Este código también está disponible en el siguiente repositorio de github:

https://github.com/pabloiglesia/ai_manager

.2.1 main.py

```
1 #!/usr/bin/env python
2 """
3 Code used to train the UR3 robot to perform a pick and place task
using Reinforcement Learning and Image Recognition.
4 This code does not perform actions directly into the robot, it
just posts actions in a ROS topic and
5 gathers state information from another ROS topic.
6 """
7
8 import rospy
9 import torch
10
11 from ai_manager.srv import GetActions, GetActionsResponse
12 from RLAlgorithm import RLAlgorithm
13 from Environment import Environment
14
15 rospy.init_node('ai_manager', anonymous=True) # ROS node
initialization
```

```
16 # Global Image Controller
17 RL_ALGORITHM = RLAlgorithm.recover_training(batch_size=256, lr
    =0.0001,
18                                     others='
    optimal_original_rewards_algorithm1901',
    include_pick_prediction=False)
19                                     # others = '
    optimal_original_rewards_new_model')
20
21 def handle_get_actions(req):
22     """
23     Callback for each Request from the Robot
24     :param req: Robot requests has 3 elements: object_gripped, x
    and y elements
25     :return:
26     """
27     object_gripped = req.object_gripped
28     current_coordinates = [req.x, req.y]
29     # Next action is calculated from the current state
30     action = RL_ALGORITHM.next_training_step(current_coordinates,
    object_gripped)
31
32     # RL_ALGORITHM.plot()
33
34     return GetActionsResponse(action)
35
36
37 def get_actions_server():
38     """
39     Service initialization to receive requests of actions from the
    robot.
40     Each time that a request is received, handle_get_actions
    function will be called
41     :return:
42     """
43     s = rospy.Service('get_actions', GetActions,
    handle_get_actions)
44     rospy.loginfo("Ready to send actions.")
45     rospy.spin()
46     rospy.on_shutdown(save_training)
47
48
49 def save_training():
50     RL_ALGORITHM.save_training()
51
52
53 if __name__ == '__main__':
54     try:
55         get_actions_server()
```

```
56     except rospy.ROSInterruptException:
57         pass
```

.2.2 RLAlgorithm.py

```
1 # coding=utf-8
2 import math
3 import random
4 import os
5 import errno
6 import sys
7 from collections import namedtuple
8
9 import matplotlib
10 import matplotlib.pyplot as plt
11 import rospy
12 import torch
13 import torch.nn as nn
14 import torch.nn.functional as F
15 import torch.optim as optim
16 import torchvision.transforms as T
17 from PIL import Image
18
19 from Environment import Environment
20 from TrainingStatistics import TrainingStatistics
21 from ImageProcessing.ImageModel import ImageModel
22 from ImageController import ImageController
23
24 is_ipython = 'inline' in matplotlib.get_backend()
25 if is_ipython: from IPython import display
26
27 import pickle
28
29 State = namedtuple( # State information namedtuple
30     'State',
31     ('coordinate_x', 'coordinate_y', 'pick_probability', '
32     object_gripped', 'image_raw')
33 )
34
35 Experience = namedtuple( # Replay Memory Experience namedtuple
36     'Experience',
37     ('state', 'coordinates', 'pick_probability', 'action', '
38     next_state', 'next_coordinates', 'next_pick_probability',
39     'reward', 'is_final_state')
40 )
41
42 class RLAlgorithm:
```

```
42     """
43     Class used to perform actions related to the RL Algorithm
44     training. It can be initialized with custom parameters or
45     with the default ones.
46
47     To perform a Deep Reinforcement Learning training, the
48     following steps have to be followed:
49
50     1. Initialize replay memory capacity.
51     2. Initialize the policy network with random weights.
52     3. Clone the policy network, and call it the target
53     network.
54     4. For each episode:
55         1. Initialize the starting state.
56         2. For each time step:
57             1. Select an action.
58                 - Via exploration or exploitation
59             2. Execute selected action in an emulator or in
60             Real-life.
61             3. Observe reward and next state.
62             4. Store experience in replay memory.
63             5. Sample random batch from replay memory.
64             6. Preprocess states from batch.
65             7. Pass batch of preprocessed states to policy
66             network.
67             8. Calculate loss between output Q-values and
68             target Q-values.
69                 - Requires a pass to the target network for
70                 the next state
71             9. Gradient descent updates weights in the policy
72             network to minimize loss.
73                 - After time steps, weights in the target
74                 network are updated to the weights in the policy network.
75
76     """
77
78     def __init__(self, object_gripped_reward=10,
79                   object_not_picked_reward=-10, out_of_limits_reward=-10,
80                   horizontal_movement_reward=-1, batch_size=32,
81                   gamma=0.999, eps_start=1, eps_end=0.01, eps_decay=0.0005,
82                   target_update=10, memory_size=100000, lr=0.001,
83                   num_episodes=1000, include_pick_prediction=False,
84                   save_training_others='optimal'):
85
86         """
87
88         :param object_gripped_reward: Object gripped reward
89         :param object_not_picked_reward: Object not picked reward
90         :param out_of_limits_reward: Out of limits reward
91         :param horizontal_movement_reward: Horizontal movement
```

```
reward
79     :param batch_size: Size of the batch used to train the
network in every step
80     :param gamma: discount factor used in the Bellman equation
81     :param eps_start: Greedy strategy epsilon start (
Probability of random choice)
82     :param eps_end: Greedy strategy minimum epsilon (
Probability of random choice)
83     :param eps_decay: Greedy strategy epsilon decay (
Probability decay of random choice)
84     :param target_update: How frequently, in terms of episodes
, target network will update the weights with the
85     policy network weights
86     :param memory_size: Capacity of the replay memory
87     :param lr: Learning rate of the Deep Learning algorithm
88     :param num_episodes: Number of episodes on training
89     :param include_pick_prediction: Use the image model pick
prediction as input of the DQN
90     :param self_training_others: Parameter used to modify the
filename of the training while saving
91     """
92
93     self.batch_size = batch_size
94     self.gamma = gamma
95     self.eps_start = eps_start
96     self.eps_end = eps_end
97     self.eps_decay = eps_decay
98     self.target_update = target_update
99     self.memory_size = memory_size
100    self.lr = lr
101    self.num_episodes = num_episodes
102    self.include_pick_prediction = include_pick_prediction
103    self.save_training_others = save_training_others
104
105    self.current_state = None # Robot current state
106    self.previous_state = None # Robot previous state
107    self.current_action = None # Robot current action
108    self.current_action_idx = None # Robot current action
Index
109    self.episode_done = False # True if the episode has just
ended
110
111    # This tells PyTorch to use a GPU if its available,
otherwise use the CPU
112    self.device = torch.device("cuda" if torch.cuda.
is_available() else "cpu") # Torch device
113    self.em = self.EnvManager(self, object_gripped_reward,
object_not_picked_reward, out_of_limits_reward,
114                                horizontal_movement_reward) #
```

```
Robot Environment Manager
115     self.strategy = self.EpsilonGreedyStrategy(self.eps_start,
116     self.eps_end, self.eps_decay) # Greedy Strategy
117     self.agent = self.Agent(self) # RL Agent
118     self.memory = self.ReplayMemory(self.memory_size) #
Replay Memory
118     self.statistics = TrainingStatistics() # Training
statistics
119
120     self.policy_net = self.DQN(self.em.image_tensor_size,
121     self.em.num_actions_available()
,
122     self.include_pick_prediction).
to(self.device) # Policy Q Network
123     self.target_net = self.DQN(self.em.image_tensor_size,
124     self.em.num_actions_available()
,
125     self.include_pick_prediction).
to(self.device) # Target Q Network
126     self.target_net.load_state_dict(self.policy_net.state_dict
()) # Target net has to be the same as policy network
127     self.target_net.eval() # Target net has to be the same as
policy network
128     self.optimizer = optim.Adam(params=self.policy_net.
parameters(), lr=self.lr) # Q Networks optimizer
129
130     print("Device: ", self.device)
131
132     class Agent:
133         """
134         Class that contains all needed methods to control the
agent through the environment and retrieve information of
135         Its state
136         """
137
138         def __init__(self, rl_algorithm):
139             """
140
141             :param self: RLAlgorithm object
142             """
143             self.strategy = rl_algorithm.strategy # Greedy
Strategy
144             self.num_actions = rl_algorithm.em.
num_actions_available() # Num of actions available
145             self.device = rl_algorithm.device # Torch device
146             self.rl_algorithm = rl_algorithm
147
148             def select_action(self, state, policy_net):
149                 """
```

```
150         Method used to pick the following action of the robot
151         Method used to pick the following action of the robot
152         :param state: State RLAlgorithm namedtuple with all
the information of the current state
153         :param policy_net: DQN object used as policy network
for the RL algorithm
154         :return:
155         """
156         random_action = False
157         if self.rl_algorithm.episode_done: # If the episode
has just ended we reset the robot environment
158             self.rl_algorithm.episode_done = False # Put the
variable episode_done back to False
159             self.rl_algorithm.statistics.new_episode()
160
161             self.rl_algorithm.current_action = 'random_state'
# Return random_state to reset the robot position
162             self.rl_algorithm.current_action_idx = None
163         else:
164             rate = self.strategy.get_exploration_rate(
165                 self.rl_algorithm.statistics.current_step) #
We get the current epsilon value
166
167             if rate > random.random(): # With a probability =
rate we choose a random action (Explore environment)
168                 action = random.randrange(self.num_actions)
169                 random_action = True
170             else: # With a probability = (1 - rate) we
Exploite the information we already have
171                 try:
172                     with torch.no_grad(): # We calculate the
action using the Policy Q Network
173                         action = policy_net(state.image_raw,
torch.tensor(
174                             [[state.coordinate_x, state.
coordinate_y]], device=self.device),
175                             state.
pick_probability, self.rl_algorithm.include_pick_prediction)\
176                             .argmax(dim=1).to(self.device) #
exploit
177                 except:
178                     print("Ha habido un error")
179
180                 self.rl_algorithm.current_action = self.
rl_algorithm.em.actions[action]
181                 self.rl_algorithm.current_action_idx = action
182
183             return self.rl_algorithm.current_action, random_action
# We return the action as a string, not as int
```

```
184
185     class DQN(nn.Module):
186         """
187         Class to create a Deep Q Learning Neural Network
188         """
189
190         def __init__(self, image_tensor_size, num_actions,
191 include_pick_prediction):
192             """
193             :param image_tensor_size: Size of the input tensor
194             :param num_actions: Number of actions, which is the
195             output of the Neural Network
196             """
197             super(RLAlgorithm.DQN, self).__init__()
198
199             self.linear1 = nn.Linear(image_tensor_size, int(
200 image_tensor_size / 2))
201             self.linear2 = nn.Linear(int(image_tensor_size / 2),
202 int(image_tensor_size / 4))
203             extra_features = 2 # coordinates
204             if include_pick_prediction:
205                 extra_features = 3 # pick prediction
206             self.linear3 = nn.Linear(int(image_tensor_size / 4) +
207 extra_features, num_actions)
208             self.linear = nn.Linear(image_tensor_size + 2,
209 num_actions)
210
211             # Called with either one element to determine next action,
212             # or a batch
213             # during optimization. Returns tensor([[left0exp,right0exp
214             ]...]).
215             def forward(self, image_raw, coordinates, pick_probability
216 =None, include_pick_probability=False):
217
218                 output = self.linear1(image_raw)
219                 output = self.linear2(output)
220                 if include_pick_probability:
221                     output = torch.cat((output, coordinates,
222 pick_probability), 1)
223                 else:
224                     output = torch.cat((output, coordinates), 1)
225                 return self.linear3(output)
226
227     class EnvManager:
228         """
229         Class used to manage the RL environment. It is used to
230         perform actions such as calculate rewards or retrieve the
231         current state of the robot.
```



```
222         """
223
224         def __init__(self, rl_algorithm, object_gripped_reward,
225 object_not_picked_reward, out_of_limits_reward,
226 horizontal_movement_reward):
227
228             """
229             Initialization of an object
230             :param rl_manager: RLAlgorithm object
231             :param object_gripped_reward: Object gripped reward
232             :param object_not_picked_reward: Object not picked
233 reward
234             :param out_of_limits_reward: Out of limits reward
235             :param horizontal_movement_reward: Horizontal movement
236 reward
237             """
238             self.object_gripped_reward = object_gripped_reward
239             self.out_of_limits_reward = out_of_limits_reward
240             self.object_not_picked_reward =
241 object_not_picked_reward
242             self.horizontal_movement_reward =
243 horizontal_movement_reward
244
245             self.device = rl_algorithm.device # Torch device
246             self.image_controller = ImageController() #
247 ImageController object to manage images
248             self.actions = ['north', 'south', 'east', 'west', '
249 pick'] # Possible actions of the objects
250             self.image_height = None # Retrieved images height
251             self.image_width = None # Retrieved Images Width
252             self.image = None # Current image ROS message
253             self.image_tensor = None # Current image tensor
254             self.pick_probability = None # Current image tensor
255
256             self.model_name = 'model-epoch=05-val_loss=0.36-
257 weights7y3_unfreeze2.ckpt'
258             # self.model_name = 'resnet50_freezed.ckpt'
259             self.model_family = 'resnet50'
260             self.image_model = ImageModel(model_name=self.
261 model_family)
262             self.feature_extraction_model = self.image_model.
263 load_model(self.model_name)
264             self.image_tensor_size = self.image_model.
265 get_size_features(
266             self.feature_extraction_model) # Size of the
267 image after performing some transformations
268
269             self.rl_algorithm = rl_algorithm
270             self.gather_image_state() # Retrieve initial state
271 image
```

```
258
259     def calculate_reward(self, previous_image):
260         """
261         Method used to calculate the reward of the previous
262         action and whether it is a final state or not
263         :return: reward, is_final_state
264         """
265         current_coordinates = [self.rl_algorithm.current_state
266                                .coordinate_x,
267                                self.rl_algorithm.current_state
268                                .coordinate_y] # Retrieve robot's current coordinates
269         object_gripped = self.rl_algorithm.current_state.
270         object_gripped # Retrieve if the robot has an object gripped
271         if Environment.is_terminal_state(current_coordinates,
272         object_gripped): # If is a terminal state
273             self.rl_algorithm.episode_done = True # Set the
274             episode_done variable to True to end up the episode
275             episode_done = True
276             if object_gripped: # If object_gripped is True,
277             the episode has ended successfully
278                 reward = self.object_gripped_reward
279                 self.rl_algorithm.statistics.
280                 add_successful_episode(True) # Saving episode successful
281                 statistic
282                 self.rl_algorithm.statistics.increment_picks()
283                 # Increase of the statistics counter
284                 rospy.loginfo("Episode ended: Object gripped!")
285             )
286             self.image_controller.record_image(
287             previous_image, True) # Saving the failure state image
288             else: # Otherwise the robot has reached the
289             limits of the environment
290                 reward = self.out_of_limits_reward
291                 self.rl_algorithm.statistics.
292                 add_successful_episode(False) # Saving episode failure
293                 statistic
294                 rospy.loginfo("Episode ended: Environment
295                 limits reached!")
296             else: # If it is not a Terminal State
297                 episode_done = False
298                 if self.rl_algorithm.current_action == 'pick': #
299                 if it is not the first action and action is pick
300                     reward = self.object_not_picked_reward
301                     self.image_controller.record_image(
302                     previous_image, False) # Saving the failure state image
303                     self.rl_algorithm.statistics.increment_picks()
304                     # Increase of the statistics counter
305                     else: # otherwise
306                         self.rl_algorithm.statistics.
```

```
fill_coordinates_matrix(current_coordinates)
288         reward = self.horizontal_movement_reward
289
290         self.rl_algorithm.statistics.add_reward(reward) # Add
reward to the algorithm statistics
291         return reward, episode_done
292
293     def gather_image_state(self):
294         """
295         This method gather the relative state of the robot by
retrieving an image using the image_controller class,
296         which reads the image from the ROS topic specified.
297         """
298         previous_image = self.image
299         self.image, self.image_width, self.image_height = self
.image_controller.get_image() # We retrieve state image
300         self.image_tensor, pick_probability = self.
extract_image_features(self.image)
301         if self.rl_algorithm.include_pick_prediction:
302             self.pick_probability = pick_probability
303
304         return previous_image
305
306     def extract_image_features(self, image):
307         """
308         Method used to transform the image to extract image
features by passing it through the image_model CNN
309         network
310         :param image_raw: Image
311         :return:
312         """
313         features, pick_prediction = self.image_model.
evaluate_image(image, self.feature_extraction_model)
314         features = torch.from_numpy(features)
315         return features.to(self.device), torch.tensor([[math.
exp(pick_prediction.numpy()[0][1]])]).to(self.device)
316
317     def num_actions_available(self):
318         """
319         Returns the number of actions available
320         :return: Number of actions available
321         """
322         return len(self.actions)
323
324     class EpsilonGreedyStrategy:
325         """
326         Class used to perform the Epsilon greede strategy
327         """
328
```

```
329         def __init__(self, start, end, decay):
330             """
331             Initialization
332             :param start: Greedy strategy epsilon start (
Probability of random choice)
333             :param end: Greedy strategy minimum epsilon (
Probability of random choice)
334             :param decay: Greedy strategy epsilon decay (
Probability decay of random choice)
335             """
336             self.start = start
337             self.end = end
338             self.decay = decay
339
340         def get_exploration_rate(self, current_step):
341             """
342             It calculates the rate depending on the actual step of
the execution
343             :param current_step: step of the training
344             :return:
345             """
346             return self.end + (self.start - self.end) * \
347                 math.exp(-1. * current_step * self.decay)
348
349         class QValues:
350             """
351             It returns the predicted q-values from the policy_net for
the specific state-action pairs that were passed in.
352             states and actions are the state-action pairs that were
sampled from replay memory.
353             """
354
355             @staticmethod
356             def get_current(policy_net, states, coordinates, actions,
pick_probabilities, include_pick_prediction = False):
357                 """
358                 With the current state of the policy network, it
calculates the q_values of
359                 :param policy_net: policy network used to decide the
actions
360                 :param states: Set of state images (Preprocessed)
361                 :param coordinates: Set of robot coordinates
362                 :param actions: Set of taken actions
363                 :return:
364                 """
365                 return policy_net(states, coordinates,
pick_probabilities, include_pick_prediction).gather(dim=1,
index=actions.unsqueeze(-1))
366
```

```
367         @staticmethod
368         def get_next(target_net, next_states, next_coordinates,
next_pick_probabilities, is_final_state,
include_pick_prediction = False):
369             """
370             Calculate the maximum q-value predicted by the
target_net among all possible next actions.
371             If the action has led to a terminal state, next reward
will be 0. If not, it is calculated using the target
372             net
373             :param target_net: Target Deep Q Network
374             :param next_states: Next states images
375             :param next_coordinates: Next states coordinates
376             :param is_final_state: Tensor indicating whether this
action has led to a final state or not.
377             :return:
378             """
379             batch_size = next_states.shape[0] # The batch size is
taken from next_states shape
380             # q_values is initialized with a zeros tensor of
batch_size and if there is GPU it is loaded to it
381             q_values = torch.zeros(batch_size).to(torch.device("
cuda" if torch.cuda.is_available() else "cpu"))
382             non_final_state_locations = (is_final_state == False)
# Non final state index locations are calculated
383             non_final_states = next_states[
non_final_state_locations] # non final state images
384             non_final_coordinates = next_coordinates[
non_final_state_locations] # non final coordinates
385             if include_pick_prediction:
386                 non_final_pick_probabilities =
next_pick_probabilities[non_final_state_locations] # non final
pick probabilities
387             else:
388                 non_final_pick_probabilities = None
389             # Max q values of the non final states are calculated
using the target net
390             q_values[non_final_state_locations] = \
391                 target_net(non_final_states, non_final_coordinates
, non_final_pick_probabilities, include_pick_prediction).max(
dim=1)[
392                 0].detach()
393             return q_values
394
395         class ReplayMemory:
396             """
397             Class used to create a Replay Memory for the RL algorithm
398             """
399
```

```
400     def __init__(self, capacity):
401         """
402         Initialization of ReplayMemory
403         :param capacity: Capacity of Replay Memory
404         """
405         self.capacity = capacity
406         self.memory = [] # Actual memory. it will be filled
with Experience namedtuples
407         self.push_count = 0 # will be used to keep track of
how many experiences have been added to the memory
408
409     def push(self, experience):
410         """
411         Method used to fill the Replay Memory with experiences
412         :param experience: Experience namedtuple
413         :return:
414         """
415         if len(self.memory) < self.capacity: # if memory is
not full, new experience is appended
416             self.memory.append(experience)
417         else: # If its full, we add a new experience and take
the oldest out
418             self.memory[self.push_count % self.capacity] =
experience
419             self.push_count += 1 # we increase the memory counter
420
421     def sample(self, batch_size):
422         """
423         Returns a random sample of experiences
424         :param batch_size: Number of randomly sampled
experiences returned
425         :return: random sample of experiences (Experience
namedtuples)
426         """
427         return random.sample(self.memory, batch_size)
428
429     def can_provide_sample(self, batch_size):
430         """
431         returns a boolean telling whether or not we can sample
from memory. Recall that the size of a sample
432         we'll obtain from memory will be equal to the batch
size we use to train our network.
433         :param batch_size: Batch size to train the network
434         :return: boolean telling whether or not we can sample
from memory
435         """
436         return len(self.memory) >= batch_size
437
438     def extract_tensors(self, experiences, include_pick_prediction
```

```
=False):
    """
439
440    Converts a batch of Experiences to Experience of batches
    and returns all the elements separately.
441    :param experiences: Batch of Experience objects
442    :return: A tuple of each element of a Experience
    namedtuple
443    """
444    batch = Experience(*zip(*experiences))
445
446    states = torch.cat(batch.state)
447    actions = torch.cat(batch.action)
448    rewards = torch.cat(batch.reward)
449    next_states = torch.cat(batch.next_state)
450    coordinates = torch.cat(batch.coordinates)
451    next_coordinates = torch.cat(batch.next_coordinates)
452    if include_pick_prediction:
453        pick_probabilities = torch.cat(batch.pick_probability)
454        next_pick_probabilities = torch.cat(batch.
next_pick_probability)
455    else:
456        pick_probabilities = None
457        next_pick_probabilities = None
458    is_final_state = torch.cat(batch.is_final_state)
459
460    return states, coordinates, pick_probabilities, actions,
rewards, next_states, next_coordinates, \
461        next_pick_probabilities, is_final_state
462
463    @staticmethod
464    def saving_name(batch_size, gamma, eps_start, eps_end,
eps_decay, lr, others=''):
465        return 'bs{}_g{}_es{}_ee{}_ed{}_lr{}_{}_pkl'.format(
466            batch_size, gamma, eps_start, eps_end, eps_decay, lr,
others
467        )
468
469    def save_training(self, dir='trainings/', others='optimal'):
470
471        filename = self.saving_name(self.batch_size, self.gamma,
self.eps_start, self.eps_end, self.eps_decay, self.lr,
472            self.save_training_others)
473
474        def create_if_not_exist(filename, dir):
475            current_path = os.path.dirname(os.path.realpath(
__file__))
476            filename = os.path.join(current_path, dir, filename)
477            if not os.path.exists(os.path.dirname(filename)):
478                try:
```

```
479         os.makedirs(os.path.dirname(filename))
480     except OSError as exc: # Guard against race
condition
481         if exc.errno != errno.EEXIST:
482             raise
483         return filename
484
485     rospy.loginfo("Saving training...")
486
487     abs_filename = create_if_not_exist(filename, dir)
488
489     self.em.image_model = None
490     self.em.feature_extraction_model = None
491
492     with open(abs_filename, 'wb+') as output: # Overwrites
any existing file.
493         pickle.dump(self, output, pickle.HIGHEST_PROTOCOL)
494
495     rospy.loginfo("Saving Statistics...")
496     print(filename)
497
498     filename = 'trainings/{}_stats.pkl'.format(filename.split(
'.pkl')[0])
499     self.statistics.save(filename=filename)
500
501     rospy.loginfo("Training saved!")
502
503     @staticmethod
504     def recover_training(batch_size=32, gamma=0.999, eps_start=1,
eps_end=0.01,
505                        eps_decay=0.0005, lr=0.001,
include_pick_prediction=False, others='optimal', dir='trainings
/', ):
506         current_path = os.path.dirname(os.path.realpath(__file__))
507         filename = RLAlgorithm.saving_name(batch_size, gamma,
eps_start, eps_end, eps_decay, lr, others)
508         filename = os.path.join(current_path, dir, filename)
509         try:
510             with open(filename, 'rb') as input:
511                 rl_algorithm = pickle.load(input)
512                 rospy.loginfo("Training recovered. Next step will
be step number {}".
513                                .format(rl_algorithm.statistics.
current_step))
514
515                 rl_algorithm.em.image_model = ImageModel(
model_name=rl_algorithm.em.model_family)
516                 rl_algorithm.em.feature_extraction_model =
rl_algorithm.em.image_model.load_model(
```



```
517         rl_algorithm.em.model_name)
518
519         return rl_algorithm
520     except IOError:
521         rospy.loginfo("There is no Training saved. New object
has been created")
522         return RLAlgorithm(batch_size=batch_size, gamma=gamma,
eps_start=eps_start, eps_end=eps_end,
523                             eps_decay=eps_decay, lr=lr,
include_pick_prediction=include_pick_prediction,
save_training_others=others)
524
525     def train_net(self):
526         """
527         Method used to train both the train and target Deep Q
Networks. We train the network minimizing the loss between
528         the current Q-values of the action-state tuples and the
target Q-values. Target Q-values are calculated using
529         thew Bellman's equation:
530
531          $q^*(state, action) = Reward + \gamma * \max(q^*(next\_state, next\_action))$ 
532         :return:
533         """
534         # If there are at least as much experiences stored as the
batch size
535         if self.memory.can_provide_sample(self.batch_size):
536             experiences = self.memory.sample(self.batch_size) #
Retrieve the experiences
537             # We split the batch of experience into different
tensors
538             states, coordinates, pick_probabilities, actions,
rewards, next_states, next_coordinates, \
539                 next_pick_probabilities, is_final_state = self.
extract_tensors(experiences, self.include_pick_prediction)
540             # To compute the loss, current_q_values and
target_q_values have to be calculated
541             current_q_values = self.QValues.get_current(self.
policy_net, states, coordinates, actions,
542                 pick_probabilities, self.include_pick_prediction)
543             # next_q_values is the maximum Q-value of each future
state
544             next_q_values = self.QValues.get_next(self.target_net,
next_states, next_coordinates,
545                 next_pick_probabilities, is_final_state, self.
include_pick_prediction)
546             target_q_values = (next_q_values * self.gamma) +
```

```
rewards
547
548         loss = F.mse_loss(current_q_values, target_q_values.
unsqueeze(1)) # Loss is calculated
549         self.optimizer.zero_grad() # set all the gradients to
0 (initialization) so that we don't accumulate
550         # gradient throughout all the backpropagation
551         loss.backward(
552             retain_graph=True) # Compute the gradient of the
loss with respect to all the weights and biases in the
553         # policy net
554         self.optimizer.step() # Updates the weights and
biases with the gradients computed
555
556         if self.statistics.episode % self.target_update == 0: #
If target_net has to be updated in this episode
557             self.target_net.load_state_dict(self.policy_net.
state_dict()) # Target net is updated
558
559     def next_training_step(self, current_coordinates,
object_gripped):
560         """
561         This method implements the Reinforcement Learning
algorithm to control the UR3 robot. As the algorithm is
prepared
562         to be executed in real life, rewards and final states
cannot be received until the action is finished, which is the
563         beginning of next loop. Therefore, during an execution of
this function, an action will be calculated and the
564         previous action, its reward and its final state will be
stored in the replay memory.
565         :param current_coordinates: Tuple of float indicating
current coordinates of the robot
566         :param object_gripped: Boolean indicating whether or not
ann object has been gripped
567         :return: action taken
568         """
569         self.statistics.new_step() # Add new steps statistics
570         self.previous_state = self.current_state # Previous state
information to store in the Replay Memory
571         previous_action = self.current_action # Previous action
to store in the Replay Memory
572         previous_action_idx = self.current_action_idx # Previous
action index to store in the Replay Memory
573         previous_image = self.em.gather_image_state() # Gathers
current state image
574
575         self.current_state = State(current_coordinates[0],
current_coordinates[1], self.em.pick_probability,
```

```
576         object_gripped, self.em.  
image_tensor) # Updates current_state  
577  
578     # Calculates previous action reward and establish whether  
the current state is terminal or not  
579     previous_reward, is_final_state = self.em.calculate_reward  
(previous_image)  
580     action, random_action = self.agent.select_action(self.  
current_state,  
581                                                     self.  
policy_net) # Calculates action  
582  
583     # There are some defined rules that the next action have  
to accomplish depending on the previous action  
584     action_ok = False  
585     while not action_ok:  
586         # Its forbidden to perform two cosecutive pick actions  
in the same place  
587         if action == 'pick' and previous_action != 'pick':  
588             action_ok = True  
589         # If previous action was south, it is forbidden to  
perform a 'north' action for  
590         # The robot not to go back to the original position.  
591         elif action == 'north' and previous_action != 'south':  
592             action_ok = True  
593         # If previous action was north, it is forbidden to  
perform a 'south' action for  
594         # The robot not to go back to the original position.  
595         elif action == 'south' and previous_action != 'north':  
596             action_ok = True  
597         # If previous action was east, it is forbidden to  
perform a 'west' action for  
598         # The robot not to go back to the original position.  
599         elif action == 'west' and previous_action != 'east':  
600             action_ok = True  
601         # If previous action was west, it is forbidden to  
perform a 'east' action for  
602         # The robot not to go back to the original position.  
603         elif action == 'east' and previous_action != 'west':  
604             action_ok = True  
605         elif action == 'random_state':  
606             action_ok = True  
607         else:  
608             action, random_action = self.agent.select_action(  
self.current_state,  
609 self.policy_net) # Calculates action  
610  
611         if random_action:
```

```
612         self.statistics.random_action() # Recollecting
statistics
613
614         # Random_state actions are used just to initialize the
environment to a random position, so it is not taken into
615         # account while storing state information in the Replay
Memory.
616         # If previous action was a random_state and it is not the
first step of the training
617         if previous_action != 'random_state' and self.statistics.
current_step > 1:
618             self.memory.push( # Pushing experience to Replay
Memory
619                             Experience( # Using an Experience namedtuple
620                                     self.previous_state.image_raw, # Initial
state image
621                                     torch.tensor([[self.previous_state.
coordinate_x, self.previous_state.coordinate_y]],
622                                                  device=self.device), # Initial
coordinates
623                                     self.previous_state.pick_probability,
624                                     torch.tensor([previous_action_idx], device=
self.device), # Action taken
625                                     self.current_state.image_raw, # Final state
image
626                                     torch.tensor([[self.current_state.coordinate_x
,
627                                                  self.current_state.coordinate_y
]],
628                                                  device=self.device), # Final
coordinates
629                                     self.current_state.pick_probability,
630                                     torch.tensor([previous_reward], device=self.
device), # Action reward
631                                     torch.tensor([is_final_state], device=self.
device) # Episode ended
632                                 ))
633
634         # Logging information
635         rospy.loginfo("Step: {}, Episode: {}, Previous reward:
{}, Previous action: {}".format(
636             self.statistics.current_step - 1,
637             self.statistics.episode,
638             previous_reward,
639             previous_action))
640
641         self.train_net() # Both policy and target networks
gets trained
642
```

```
643         return action
```

.2.3 Environment.py

```
1  """
2  This class defines a RL environment for a pick and place task with
   a UR3 robot.
3  This environment is defined by its center (both cartesian and
   angular coordinates), the total length of its x and y axis
4  and other parameters
5  """
6
7  import random
8  from BlobDetector.BlobDetector import BlobDetector
9  from ai_manager.ImageController import ImageController
10 from math import floor
11
12
13 class Environment:
14     X_LENGTH = 0.175 # Total length of the x axis environment in
   meters
15     Y_LENGTH = 0.225 # Total length of the y axis environment in
   meters
16
17     CAMERA_SECURITY_MARGIN = 0.035 # As the camera is really
   close to the gripping point, it needs a security marging
18     X_LIMIT = X_LENGTH - CAMERA_SECURITY_MARGIN # Robot
   boundaries of movement in axis X
19     Y_LIMIT = Y_LENGTH - CAMERA_SECURITY_MARGIN # Robot
   boundaries of movement in axis Y
20
21     CARTESIAN_CENTER = [-0.31899288568, -0.00357907370787,
   0.376611799631] # Cartesian center of the RL environment
22     ANGULAR_CENTER = [2.7776150703430176, -1.5684941450702112,
   1.299912452697754, -1.3755658308612269,
23     -1.5422008673297327, -0.3250663916217249] #
   Angular center of the RL environment
24     PLACE_CARTESIAN_CENTER = [0, 0.25, CARTESIAN_CENTER[2]] #
   Cartesian center of the place box
25     ANGULAR_PICTURE_PLACE = [1.615200161933899,
   -1.235102955495016, 0.739865779876709, -1.2438910643206995,
   -1.5095704237567347, -0.06187755266298467]
26
27     PICK_DISTANCE = 0.01 # Distance to the object when the robot
   is performing the pick and place action
28     ACTION_DISTANCE = 0.02 # Distance to the object when the
   robot is performing the pick and place action
29
```

```
30     ENV_BOUNDS_TOLERANCE = 0
31
32     @staticmethod
33     def generate_random_state(image=None, strategy='ncc'):
34         """
35         Calculates random coordinates inside the Relative
36         Environment defined.
37         To help the robot empty the box, the generated coordinates
38         won't be in the center of the box, because this is
39         the most reachable place of the box.
40
41         :param strategy: strategy used to calculate random_state
42         coordinates
43         :return:
44         """
45         def generate_random_coordinates():
46             coordinate_x = random.uniform((-Environment.X_LIMIT +
47             Environment.ENV_BOUNDS_TOLERANCE) / 2,
48             (Environment.X_LIMIT -
49             Environment.ENV_BOUNDS_TOLERANCE) / 2)
50             coordinate_y = random.uniform((-Environment.Y_LIMIT +
51             Environment.ENV_BOUNDS_TOLERANCE) / 2,
52             (Environment.Y_LIMIT -
53             Environment.ENV_BOUNDS_TOLERANCE) / 2)
54             return coordinate_x, coordinate_y
55
56         # Random coordinates avoiding the ones in the center,
57         # which have a bigger probability of being reached by the
58         # robot.
59         if strategy == 'ncc' or strategy == '
60         non_centered_coordinates':
61             coordinates_in_center = True
62             while coordinates_in_center:
63                 coordinate_x, coordinate_y =
64                 generate_random_coordinates()
65                 if abs(coordinate_x) > (Environment.X_LIMIT / 4)
66                 or abs(coordinate_y) > (Environment.Y_LIMIT / 4):
67                     coordinates_in_center = False
68             elif strategy == 'optimal' and image is not None: #
69                 Before going to a random state, we check that there are pieces
70                 in this place
71                 blob_detector = BlobDetector(x_length=Environment.
72                 X_LENGTH, y_length=Environment.Y_LENGTH, columns=4, rows=4)
73                 optimal_quadrant = blob_detector.find_optimal_quadrant
74                 (image)
75                 optimal_point = blob_detector.quadrants_center[
76                 optimal_quadrant]
77
78                 coordinate_x = optimal_point[0] * 0.056
```

```
63         coordinate_y = optimal_point[1] * 0.056
64     else: # Totally random coordinates
65         coordinate_x, coordinate_y =
generate_random_coordinates()
66
67     return [coordinate_x, coordinate_y]
68
69     @staticmethod
70     def get_relative_corner(corner):
71         """
72         Function used to calculate the coordinates of the
environment corners relative to the CARTESIAN_CENTER.
73
74         :param corner: it indicates the corner that we want to get
the coordinates. It's composed by two letters
75         that indicate the cardinality. For example: ne indicates
North-East corner
76         :return coordinate_x, coordinate_y:
77         """
78         if corner == 'sw' or corner == 'ws':
79             return -Environment.X_LIMIT / 2, Environment.Y_LIMIT /
2
80         if corner == 'nw' or corner == 'wn':
81             return Environment.X_LIMIT / 2, Environment.Y_LIMIT /
2
82         if corner == 'ne' or corner == 'en':
83             return Environment.X_LIMIT / 2, -Environment.Y_LIMIT /
2
84         if corner == 'se' or corner == 'es':
85             return -Environment.X_LIMIT / 2, -Environment.Y_LIMIT
/ 2
86
87     @staticmethod
88     def is_terminal_state(coordinates, object_gripped):
89         """
90         Function used to determine if the current state of the
robot is terminal or not
91         :return: bool
92         """
93         def get_limits(length): return length / 2 - Environment.
ENV_BOUNDS_TOLERANCE # funcion to calculate the box boundaries
94         x_limit_reached = abs(coordinates[0]) > get_limits(
Environment.X_LIMIT) # x boundary reached
95         y_limit_reached = abs(coordinates[1]) > get_limits(
Environment.Y_LIMIT) # y boundary reached
96         return x_limit_reached or y_limit_reached or
object_gripped # If one or both or the boundaries are reached
--> terminal state
```

.2.4 ImageController.py

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 import os
5 import time
6
7 import rospy
8 from PIL import Image as PILImage
9 from sensor_msgs.msg import Image
10
11 """
12 This class is used to manage sensor_msgs Images.
13 """
14
15 class ImageController:
16     def __init__(self, path=os.path.dirname(os.path.realpath(
17         __file__)), image_topic='/usb_cam/image_raw'):
18         self.ind_saved_images = 0 # Index which will tell us the
19         number of images that have been saved
20         self.success_path = "{}/{}/success".format(path) # Path
21         where the images are going to be saved
22         self.fail_path = "{}/{}/fail".format(path) # Path where the
23         images are going to be saved
24         self.image_topic = image_topic
25
26         # If it does not exist, we create the path folder in our
27         workspace
28         try:
29             os.stat(self.success_path)
30         except:
31             os.mkdir(self.success_path)
32
33         # If it does not exist, we create the path folder in our
34         workspace
35         try:
36             os.stat(self.fail_path)
37         except:
38             os.mkdir(self.fail_path)
39
40     def get_image(self):
41         msg = rospy.wait_for_message(self.image_topic, Image)
42
43         return self.to_pil(msg), msg.width, msg.height
44
45     def record_image(self, img, success):
46         path = self.success_path if success else self.fail_path #
47         The path were we want to save the image is
```



```
41
42     image_path = '{} /img{}.png'.format( # Saving image
43         path, # Path
44         time.time()) # FIFO queue
45
46     img.save(image_path)
47
48     self.ind_saved_images += 1 # Index increment
49
50     def to_pil(self, msg, display=False):
51         size = (msg.width, msg.height) # Image size
52         img = PILImage.frombytes('RGB', size, msg.data) #
53         sensor_msg to Image
54         return img
55
56 if __name__ == '__main__':
57     rospy.init_node('image_recorder') # ROS node initialization
58     image_controller = ImageController(path='/home/pilar/Pilar/
59     ros_pictures', image_topic='/usb_cam2/image_raw')
60     img, width, height = image_controller.get_image()
61     image_controller.record_image(img, True)
```

BIBLIOGRAPHY

- [1] *Formación en línea de CB3*. URL: <https://academy.universal-robots.com/es/formacion-en-linea-gratuita/formacion-en-linea-de-cb3/> (visited on 11/15/2020).
- [2] Preferred Networks, Inc. *Bin-picking Robot Deep Learning*. 2015. URL: https://www.youtube.com/watch?v=ydh_AdWZflA&ab_channel=Pickit3D (visited on 11/27/2020).
- [3] Guillaume Lample and Devendra Singh Chaplot. “Playing FPS Games with Deep Reinforcement Learning”. In: *arXiv:1609.05521 [cs]* (Jan. 29, 2018). arXiv: [1609.05521](https://arxiv.org/abs/1609.05521). URL: <http://arxiv.org/abs/1609.05521> (visited on 11/30/2020).
- [4] Y. Zhu et al. “Target-driven visual navigation in indoor scenes using deep reinforcement learning”. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. 2017 IEEE International Conference on Robotics and Automation (ICRA). May 2017, pp. 3357–3364. DOI: [10.1109/ICRA.2017.7989381](https://doi.org/10.1109/ICRA.2017.7989381).
- [5] *Reinforcement Learning - Goal Oriented Intelligence*. URL: https://deeplizard.com/learn/playlist/PLZbbT5o_s2xoWNVdDudn51XM8lOuZ_Njv (visited on 11/28/2020).
- [6] OpenAI. *Gym: A toolkit for developing and comparing reinforcement learning algorithms*. URL: <https://gym.openai.com> (visited on 11/30/2020).
- [7] A. Rupam Mahmood et al. “Setting up a Reinforcement Learning Task with a Real-World Robot”. In: *arXiv:1803.07067 [cs, stat]* (Mar. 19, 2018). arXiv: [1803.07067](https://arxiv.org/abs/1803.07067). URL: <http://arxiv.org/abs/1803.07067> (visited on 11/30/2020).
- [8] Xiaoyun Lei, Zhian Zhang, and Peifang Dong. “Dynamic Path Planning of Unknown Environment Based on Deep Reinforcement Learning”. In: *Journal of Robotics* 2018 (Sept. 18, 2018), pp. 1–10. ISSN: 1687-9600, 1687-9619. DOI: [10.1155/2018/5781591](https://doi.org/10.1155/2018/5781591). URL: <https://www.hindawi.com/journals/jr/2018/5781591/> (visited on 11/30/2020).

- [9] Marco Wiering. “Reinforcement Learning in Dynamic Environments using Instantiated Information”. In: (Aug. 27, 2001).
- [10] *Agile Methodology: What is Agile Software Development Model?* URL: <https://www.guru99.com/agile-scrum-extreme-testing.html> (visited on 12/02/2020).