

UDACITY

PROJECT REPORT

---

## P2 - Continuous Control

---

*Author:*

Thomas DI MARTINO

*Supervisor:*

Udacity

*A report submitted in fulfilment of the requirements  
for the Student in the Deep Reinforcement Learning Nanodegree Program  
of the*

Artificial Intelligence department

May 2020



# Declaration of Authorship

I, Thomas DI MARTINO, confirm that this work, 'P2 - Continuous Control', submitted for assessment is my own and is expressed in my own words. Any uses made within it of the works of other authors in any form (e.g., ideas, equations, figures, text, tables, programs) are properly acknowledged at any point of their use. A list of the references employed is included.

Signed:

---

Date:

---

# Contents

<b>Declaration of Authorship</b>	<b>i</b>
<b>Contents</b>	<b>ii</b>
<b>List of Figures</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem diagnosis . . . . .	1
1.2 Environment Introduction, a RL perspective . . . . .	2
1.2.1 States . . . . .	2
1.2.2 Rewards . . . . .	2
1.2.3 Actions . . . . .	2
1.2.4 Success Criteria . . . . .	2
1.3 Report Structure . . . . .	3
1.4 Code location . . . . .	3
<b>2 Models &amp; Results</b>	<b>4</b>
2.1 Retained models . . . . .	4
2.1.1 DDPG Agent . . . . .	4
2.1.1.1 Hyper-Parameters and formula presentation . . . . .	4
2.1.1.2 The Critic network . . . . .	5
2.1.1.3 The Actor network . . . . .	6
2.1.1.4 Deep Learning Model Architecture . . . . .	7
2.1.1.5 File architecture . . . . .	8
2.2 Results . . . . .	8
<b>3 Conclusion</b>	<b>10</b>
3.1 Summary . . . . .	10
3.2 Overture . . . . .	10
<b>Bibliography</b>	<b>11</b>

# List of Figures

1.1	Screenshot of the agent's environment . . . . .	1
1.2	Agent state example . . . . .	2
2.1	Schematic representation of a critic network . . . . .	6
2.2	Schematic representation of an actor network . . . . .	7
2.3	Critic Model Architecture . . . . .	7
2.4	Actor Model Architecture . . . . .	8
2.5	DDPG architecture training performance with uniform replay buffer . . .	9
2.6	Prioritised DDPG architecture training performance . . . . .	9

# Chapter 1

## Introduction

### 1.1 Problem diagnosis

In the context of the *Reacher* environment, our agent is an arm with two joints that needs to move to different target locations, modeled by a round 3D Green ball rotating around the arm's extremity. The Reacher environment of the agent is highly similar to the one developed by Unity in the following sets of learning environment: <https://github.com/Unity-Technologies/ml-agents/blob/master/docs/Learning-Environment-Examples.md#reacher>

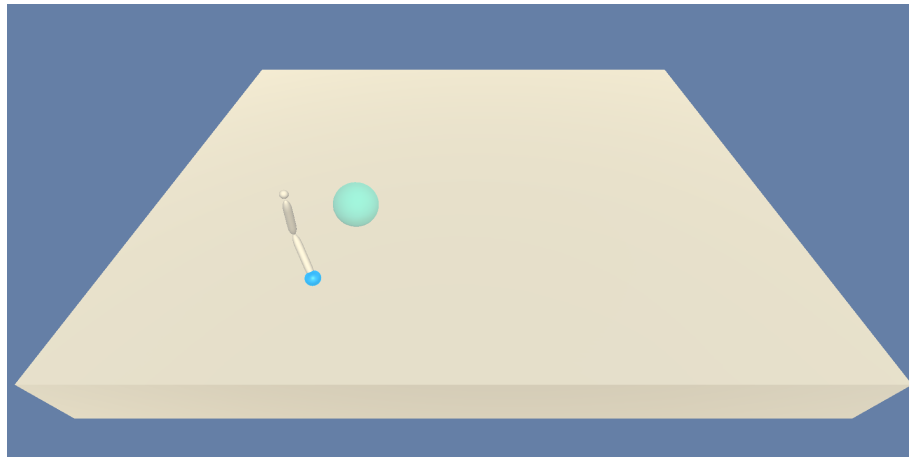


FIGURE 1.1: Screenshot of the agent's environment

## 1.2 Environment Introduction, a RL perspective

### 1.2.1 States

To capture is knowledge of the environment, our agent has an observation space of size 33. Amongst all the variables making up this vector, we can count the position of the arm, its velocity (of the arm as a whole and of its angles), its rotation index and so on.

---

```
[ 0.00000000e+00 -4.00000000e+00  0.00000000e+00  1.00000000e+00
-0.00000000e+00 -0.00000000e+00 -4.37113883e-08  0.00000000e+00
 0.00000000e+00  0.00000000e+00  0.00000000e+00  0.00000000e+00
 0.00000000e+00  0.00000000e+00 -1.00000000e+01  0.00000000e+00
 1.00000000e+00 -0.00000000e+00 -0.00000000e+00 -4.37113883e-08
 0.00000000e+00  0.00000000e+00  0.00000000e+00  0.00000000e+00
 0.00000000e+00  0.00000000e+00  5.75471878e+00 -1.00000000e+00
 5.55726671e+00  0.00000000e+00  1.00000000e+00  0.00000000e+00
-1.68164849e-01]
```

---

FIGURE 1.2: Agent state example

### 1.2.2 Rewards

The reward of the agent is pretty straight forward:

it is given +0.1 score for every time-step where the agent's hand is inside the ball. As no negative reward is given to the agent, its behaviour may not be physically-speaking ideal:

any torque will not be penalised. However, such bad behaviour should always be taken into account when developing a RL agent for a concrete robot implementation. A more elaborate reward function could then be helpful to train a better agent, even if training would then be made harder.

### 1.2.3 Actions

To collect these rewards, our agent can move itself. These movement are represented as a vector of 4 values, with each value being between -1 and 1.

### 1.2.4 Success Criteria

When training the agent, it is considered successful to the eyes of the assignment when it gathers a reward of **+30** over 100 consecutive episodes.

## 1.3 Report Structure

In this report, we will present our implementation as well as expose their resulting learning performance. We will then explore eventual amelioration of the current work, exposing the weaknesses of both the retained model as well as for the environment (a parallel may even be drawn with real life situations).

## 1.4 Code location

Located in a Github repository at the following address [https://github.com/dimartinot/P2-Continuous\\_Control](https://github.com/dimartinot/P2-Continuous_Control)

## Chapter 2

# Models & Results

### 2.1 Retained models

One agent was trained during this work:

- A DDPG <sup>1</sup> Agent that implements **prioritised experience replay** using the following formula  $\frac{p^\alpha}{\sum_i p_i^\alpha}$  where p is the lastly measured td error of the experience tuple, before its insertion in the buffer. Alpha is an hyper-parameter: the closer alpha is to zero, the more uniform the sampling distribution will be. The closer to one, the less uniform it will be.

#### 2.1.1 DDPG Agent

##### 2.1.1.1 Hyper-Parameters and formula presentation

Member of the family of Actor-Critic methods, the DDPG Agent is ideal for continuous task. Actor-Critic methods are a family of RL algorithms mixing policy-based and value-based learning:

- In **value-based** methods, the agent is learning to estimate the value of actions and situations;
- In **policy-based** methods, the agent is learning to act.

Hence, Actor-Critic methods try to learn both of these concepts, through the use of two separate neural networks:

---

<sup>1</sup>Deep Deterministic Policy Gradient



1. The **Actor**, outputting the probability of taking any action given the state and the policy;
2. The **Critic**, giving the value of an input state.

Actor-Critic methods follow this algorithm:

- Step 1: input current state  $S$  into the actor, get the action  $a$  to perform (stochastically);
- Step 2: observe next state  $S'$  and reward  $r$  - collect a tuple  $(S, a, r, S')$ ;
- Step 3: Train the critic with  $r + \gamma * V(S')$ ;
- Step 4: Use the critic to calculate the advantage  $A(S, a) = r + \gamma * V(S'; \theta_v) - V(S; \theta_v)$
- Step 5: Train the actor with the calculated advantage  $A(S, a)$ .

From this algorithm, the DDPG implementation varies in the following points:

- the actor **always** output the best possible action (no stochastic choice of action);
- the critic evaluates the actor's best-believed action;
- the DDPG algorithm uses **soft updates** and mix in 0.1% of the regular weights in the target weights at **every time step**.

Hence, given that the DDPG uses 2 networks (actor and critic networks) and that each of them as their mutual *regular* and *target* version, a DDPG algorithm manipulates a total of 4 distinct neural networks.

**The Critic network** The critic network was trained with the following hyper-parameters:

- Adam optimiser with a Learning Rate of  $10^{-3}$ ;
- A batch size of 128;
- Soft update rate  $\tau = 10^{-3}$ ;
- Delayed update: the networks are updated 20 times every 10 time step;

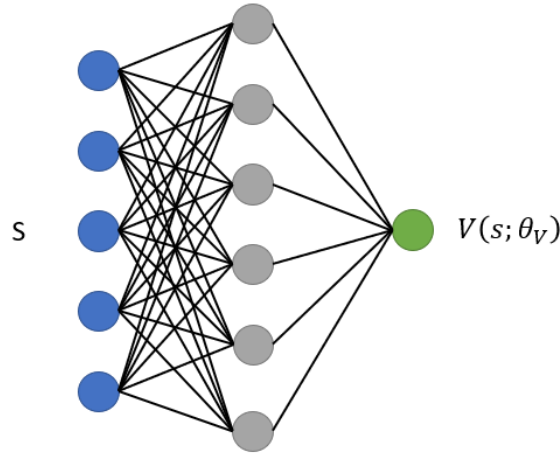


FIGURE 2.1: Schematic representation of a critic network

- Epsilon value for noise addition is set at  $\epsilon = 1.0$  at the start of training and is slowly decreasing by the subtraction of  $10^{-6}$  at every timestep;
- MSE Loss between the TD target and the local weights evaluation of the (state, action) tuple. Formally, it can be written as:  $L = \frac{1}{N} \sum_{n=1}^N (\textcolor{brown}{R}_{t+1} + \gamma * \textcolor{brown}{max}_a(\hat{q}(S_{n+1}, a, W_t)) - \textcolor{violet}{\hat{q}}(S_n, A, W_l))^2$

where:

- $\textcolor{brown}{R}_{t+1} + \gamma * \textcolor{brown}{max}_a(\hat{q}(S_{n+1}, a, W_t))$  is the target value (which is evolving through computation), that we would write as  $\hat{y}$  in a common Machine Learning problem;
- $\textcolor{violet}{\hat{q}}(S_n, A, W_l)$  is the predicted value, that we would write as  $y$ ;
- N is the batch size (in our implementation, we selected 128 as a batch size).

**The Actor network** The actor network was trained with the following hyperparameters:

- Adam optimiser with a Learning Rate of  $10^{-3}$ ;
- A batch size of 128;
- Soft update rate  $\tau = 10^{-3}$ ;
- Delayed update: the networks are updated 20 times every 10 time step;
- Epsilon value for noise addition is set at  $\epsilon = 1.0$  at the start of training and is slowly decreasing by the subtraction of  $10^{-6}$  at every timestep;

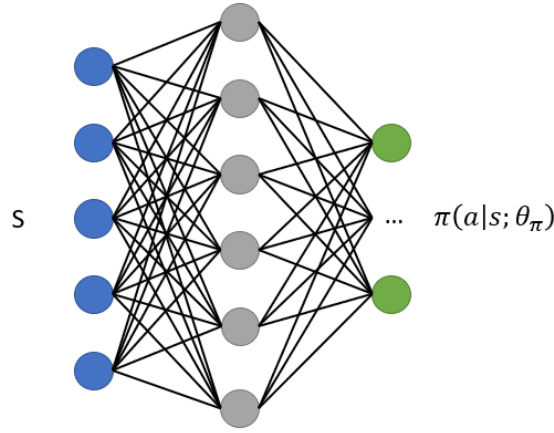


FIGURE 2.2: Schematic representation of an actor network

- The evaluation by the critic of the action chosen by the actor acts as a loss: set negative, we try to minimise it (transformed as a negative function, minimising it is the same as maximising the positive function).

To train the local network, we use a *Mean Squared Error* loss between the TD Target and the local weights evaluation of the (state, action) tuple.

### 2.1.1.2 Deep Learning Model Architecture

The following figures (cf. 2.3 & 2.4) presents the retained Deep Learning architectures. I took massive inspiration from the model used in the lab of ddpq. Experimentation was done by adding/removing layers to have a good balance between depth and shallowness.



FIGURE 2.3: Critic Model Architecture



FIGURE 2.4: Actor Model Architecture

### 2.1.1.3 File architecture

Relevant files for this part of the implementation are:

- *ddpg\_agent.py*: Contains the class definition of the basic DPPG learning algorithm that uses soft update (for weight transfer between the local and target networks) as well as a uniformly distributed replay buffer and OUNoise to model the exploration/exploitation dilemma;
- *model.py*: Contains the PyTorch class definition of the Actor and the critic neural networks, used by their mutual target and local network's version.
- *prioritized\_ddpg\_agent.py*: Contains the class definition of the DDPG learning algorithm with prioritized replay buffer.

## 2.2 Results

We now present the results of training these two architecture with the curve of the *reward per episode* function. In the following figure (cf. 2.5), two graphs are plot:

- The blue graph represents the score of the agent at the end of each episode.
- The orange graph represents the moving average score of the agent. It better represents the overall learning behaviour of the agent.

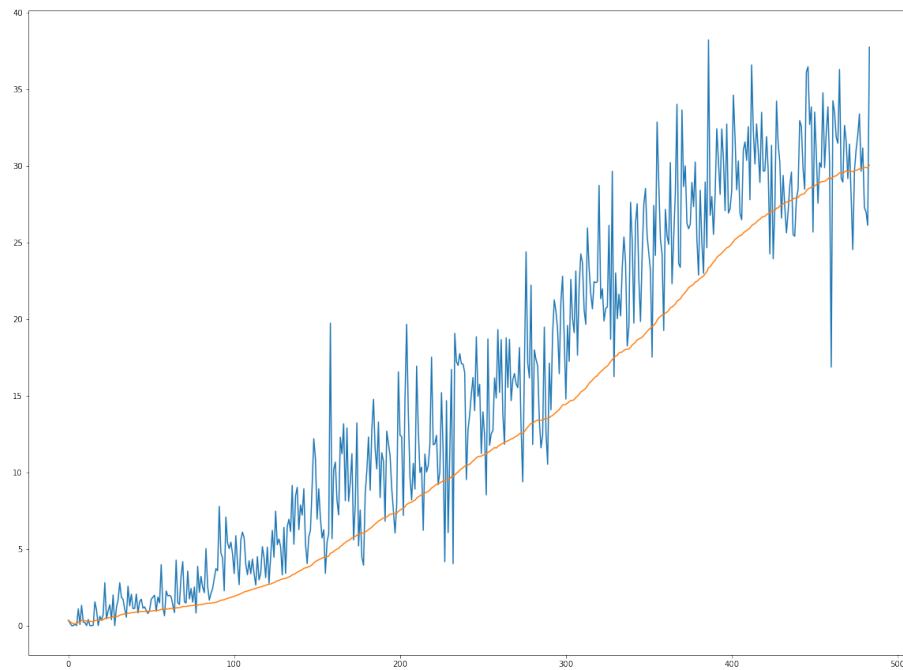


FIGURE 2.5: DDPG architecture training performance with uniform replay buffer

Our learning was pretty stable throughout the whole process, with occasional drop from time to time. However, the moving average perfectly describes the learning tendency followed by our agent.

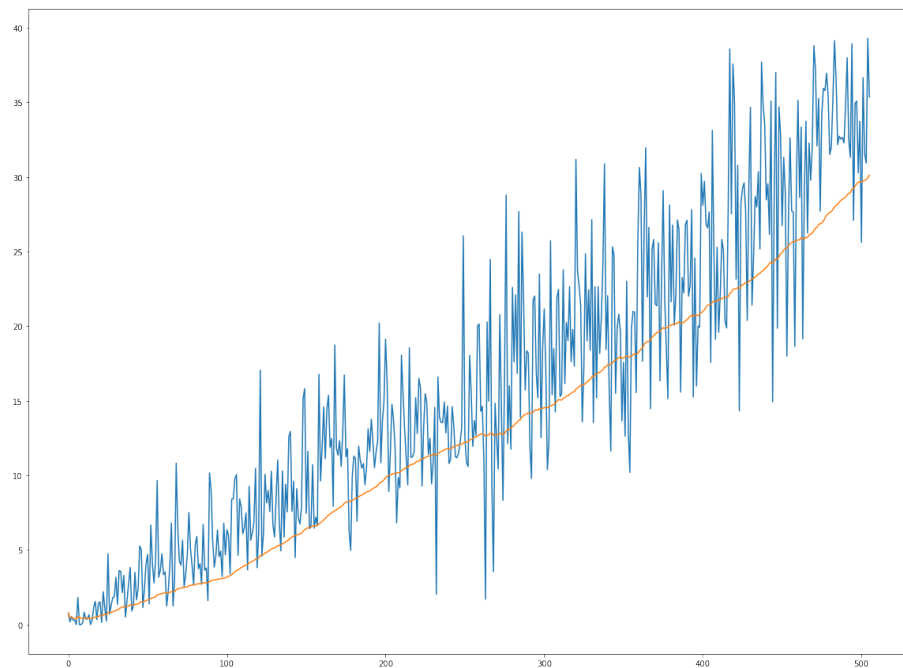


FIGURE 2.6: Prioritised DDPG architecture training performance

## Chapter 3

# Conclusion

### 3.1 Summary

As a summary, two methods were studied and used in this project: while sensibly similar performances were obtained, no thorough conclusion can be drawn with regards to their relative performance. Indeed, as not a lot of hyper parameter exploration has been done, we can only base ourselves on these limited results. Furthermore, the threshold of  $+13$  score to be obtained limited our analysis of how far can both model go. However, this was **not** the goal of this project and every objectives were met.

### 3.2 Overture

An interesting concept to be explored is the feature extraction phase with the use of a state represented as a picture of the environment. Using convolutional layers, we would see an harder task unfolds before us, which would then be very interesting to compare to the other format of this project. The main guesses in the comparison are:

- The use of the 37-sized state vector would lead to faster and better converging algorithm (as well as easier to train);
- The use of the image input would lead to better overall performance, to the cost of an harder training pipeline as well as more computation involved.

However, both of these ideas are assumptions and are yet to have been studied by myself.

# Bibliography