

# Report: Project2

## Deep Reinforced Learning Nanodegree: Continuous Control with DDPG

### Introduction

This project solves for actuation of a 2 degree of freedom robotic arm to have the end actuator reach the goal position. Implemented using DDPG (Deep Deterministic Policy Gradient) algorithm. This is a model free learning method, is pretty cool and easily portable across different applications.

### Environment

Reacher environment from Unity Machine Learning Agent Toolkit is used for this project. In this environment, a double-jointed arm can move to target locations. A reward of +0.1 is provided for each step that the agent's hand is in the goal location. Thus, the goal of your agent is to maintain its position at the target location for as many time steps as possible.

The observation space consists of 33 variables corresponding to position, rotation, velocity, and angular velocities of the arm. Each action is a vector with four numbers, corresponding to torque applicable to two joints. Every entry in the action vector is a number between -1 and 1.

State/ Observation size: 33

Action space size : 4

Goal: The environment is considered solved, when the average (over 100 episodes) of those average scores (all 20 agents) is at least +30.

### Learning Algorithm

DDPG - Deep Deterministic Policy Gradient

---

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .  
Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$   
Initialize replay buffer  $R$   
**for** episode = 1, M **do**  
    Initialize a random process  $\mathcal{N}$  for action exploration  
    Receive initial observation state  $s_1$   
    **for** t = 1, T **do**  
        Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise  
        Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$   
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$   
        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$   
        Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'}))|\theta^{Q'}$   
        Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$   
        Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

Update the target networks:

$$\begin{aligned}\theta^{Q'} &\leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \\ \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}\end{aligned}$$

**end for**  
**end for**

---

DDPG is the an Actor-Critic algorithm that simultaneously learns a policy and an action-value function  $Q(s,a)$ . The actor takes a deterministic action following a policy and critic critics the actors action based on its action-value function.

Noise is added to the action intentionally to match up for a continuous action space use case through the output is deterministic.

## Hyperparameters

<code>GAMMA</code>	<code>0.99</code>	Discount factor
<code>TAU</code>	<code>1e-3</code>	Soft update of target parameters
<code>LR_ACTOR</code>	<code>1e-3</code>	Learning rate of the actor
<code>LR_CRITIC</code>	<code>1e-3</code>	Learning rate of the actor
<code>WEIGHT_DECAY</code>	<code>0.0000</code>	L2 weight decay
<code>BATCH_SIZE</code>	<code>512</code>	Batch size
<code>BUFFER_SIZE</code>	<code>int(1e6)</code>	Replay buffer size
<code>learn_every</code>	<code>20</code>	How often to update target network
<code>num_learn</code>	<code>10</code>	Number of updates at each step
<code>OUnoise.theta</code>	<code>0.15</code>	Ornstein-Uhlenbeck Noise parameter
<code>OUnoise.mu</code>	<code>0.2</code>	Ornstein-Uhlenbeck Noise parameter

## Network Structure

### Actor : NN

a\_FC1 : state\_size - > 200

a\_FC1 : ReLU ( Batch Normalization ( a\_FC1 ) )

a\_FC2 : ReLU ( 200 (a\_FC1) -> 100 )

a\_FC3 : tanh ( 100 (a\_FC2) -> action\_size )

### Critic: NN

c\_FC1 : state\_size - > 200

c\_FC1 : ReLU ( Batch Normalization ( c\_FC1 ) )

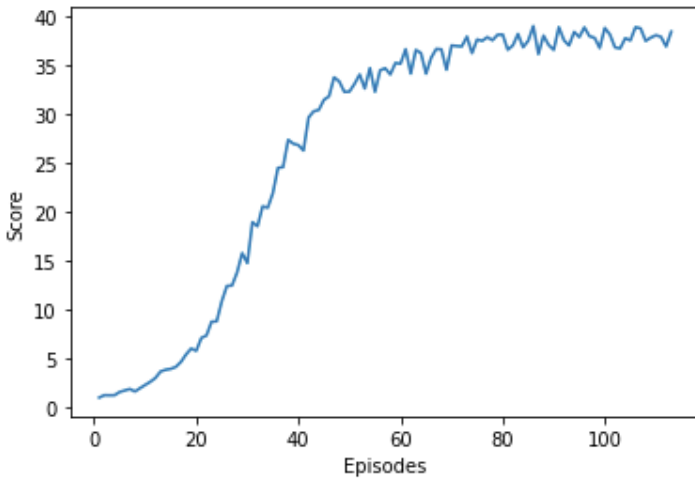
c\_FC2 : ReLU ( 200 ([ c\_FC1 <concat> a\_FC3 ]) -> 100)

c\_FC3 : 100 (c\_FC2) -> 1

## Result :

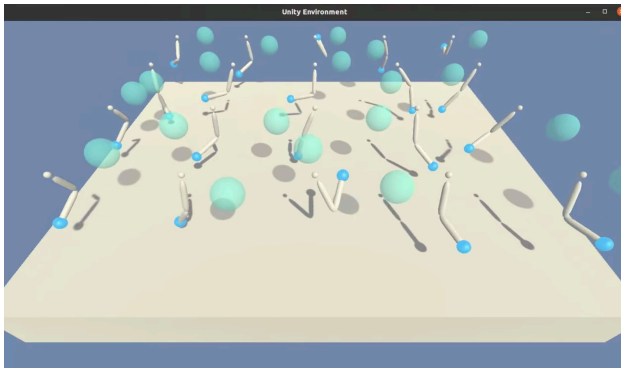
Goal was achieved after 113 episodes with total average score of 30.05

## Score Vs Episode Plot

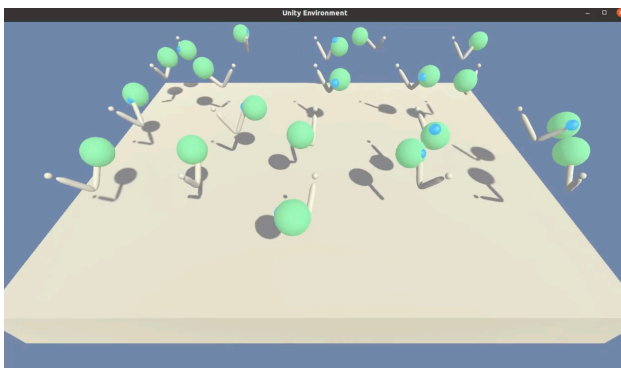


## Result GIF:

Untrained: <https://gph.is/g/Z866Pbb>



Trained: <https://gph.is/g/aRVV5qw>



## **Future Work :**

- Try out the crawler environment
- Use prioritized experience replay
- Use tensorboard to compare results form multiple hyper parameter choices