

# Deep Deterministic Policy Gradients (DDPG)

## Continuous Control Project

Algorithm Pseudocode - Source: <https://arxiv.org/pdf/1509.02971.pdf>

---

**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**

---

## What DDPG is about:

**DDPG can be thought as the extension of the Q-Learning algorithm for continuous tasks where discrete state/actions may not suffice.**

The actor will approximate a deterministic policy, in contrast to other actor networks we may come across in actor-critic methods, thus the actor will always map the “best action” whenever we query the network.

The critic is used to approximate the maximizer over the Q values of the next state (SarsaMax – Q-Learning) and not as a learned baseline, as with other actor-critic algorithms.

DDPG implements a replay buffer, much like in DQNs where the SARS experiences are saved.

## About the Project

This project uses the **Deep Deterministic Policy Gradients (DDPG)** algorithm to utilize an actor-critic approach to solve the environment using both the 1 agent and the 20 agents approaches.

## Project Files Involved

There are 4 total files involved for this project

- **model.py** contains the models for the actor and the critic neural networks
- **ddpg\_agent.py** contains the code defining the Actor class, accompanied by the ReplayBuffer and the OUNoise classes.
- **Continuous\_Control\_1\_agent.ipynb** contains the notebook with the solution for the 1 Agent Environment approach
- **Continuous\_Control\_20\_agents.ipynb** contains the notebook with the solution for the 20 Agents Environment approach

As you can tell by the files, there are two notebooks that utilize the same architecture for both approaches. The notebooks are intentionally identical, so they can display how more difficult and unstable the 1 actor approach is compared to the 20 actors approach.

The only difference between the two notebooks is the line that calls the environment, as well as the lines that save and load the checkpoints (*so we have different neural network checkpoints for each approach*).

```
UNITY_PATH_01_AGENTS = 'Reacher_01_Windows_x86_64/Reacher.exe'
```

```
UNITY_PATH_20_AGENTS = 'Reacher_20_Windows_x86_64/Reacher.exe'
```

The line that loads the environment at **Continuous\_Control\_1\_agent.ipynb** file

```
env = UnityEnvironment(file_name = UNITY_PATH_01_AGENTS)
```

The line that loads the environment at **Continuous\_Control\_20\_agents.ipynb** file

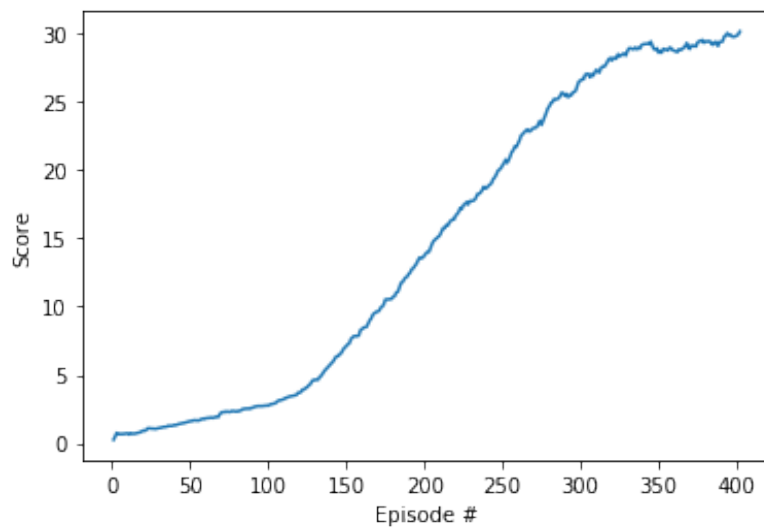
```
env = UnityEnvironment(file_name = UNITY_PATH_01_AGENTS)
```

## Comparing the training using 1 and 20 Agents

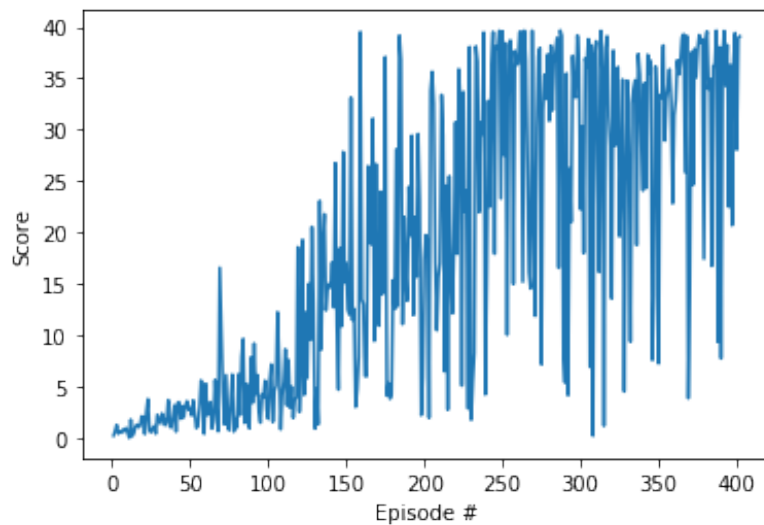
Using 1 agent, the problem was solved in 402 episodes, compared to the 20 agents approach, where **the same architecture** solved the environment in 114 episodes!

### 1 Agent approach: 402 episodes

*Agents average score over 100 episodes*

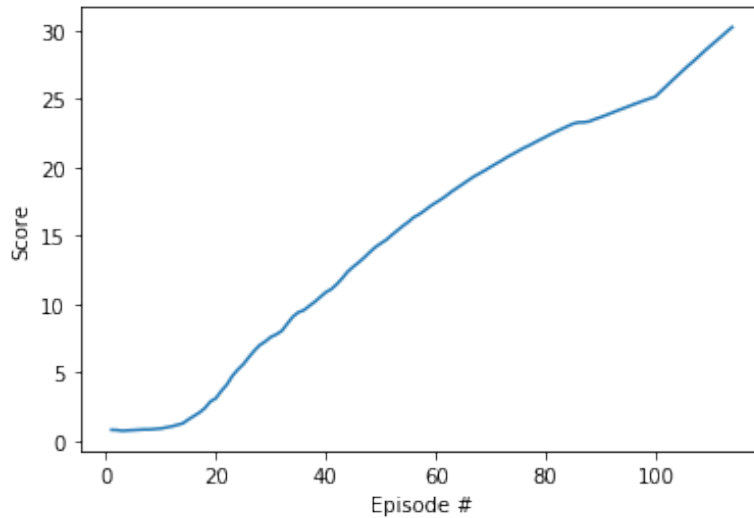


*Agents average score per episode*

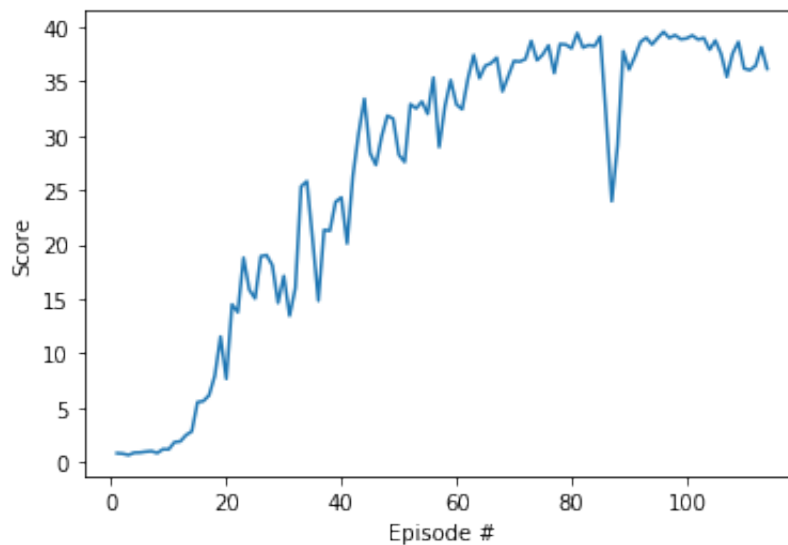


## 20 Agents approach: 114 episodes

*Agents average score over 100 episodes*



*Agents average score per episode*



It is very clear that the approach using 20 agents is way more stable, efficient, and converges much faster.

20 Agents approach takes the advantage of utilizing 20 different samples per network weights, allowing it to generalize much better by mimicking the **experience replay** functionality.

## The Architecture

There have been two neural networks for the utilization of the DDPG algorithm.

I defined the **actor network** using three internal fully connected layers of dept 512, 256, 128 depths respectively, with batch normalization to their inputs.

The **critic network** utilized two internal layers of 512 and 256 layers, with batch normalization to the inputs as well.

Another very important factor was the **noise** at the “OUNoise” function. Starting with a noise of 0.1, the system was too unstable. Started to gradually reduce the noise, where at value 0.02 it seemed that the network behaved very efficiently.

Gradient Clipping appeared to be much beneficial to the critic network, as shown in the lines bellow

```
self.critic_optimizer.zero_grad()
critic_loss.backward()
torch.nn.utils.clip_grad_norm(self.critic_local.parameters(), 1)
self.critic_optimizer.step()
```

The training after around 100 episodes was becoming very unstable, especially on the 1 agent approach **Updating the networks 10 times after every 20 timesteps** helped stabilizing the training for much longer.

## Hyper Params

```
BUFFER_SIZE = int(1e6) # replay buffer size
BATCH_SIZE = 256       # minibatch size
GAMMA = 0.9            # discount factor
TAU = 1e-3             # for soft update of target parameters
LR_ACTOR = 1e-3        # learning rate of the actor
LR_CRITIC = 1e-3       # learning rate of the critic
WEIGHT_DECAY = 0       # L2 weight decay
```

## What's next

It seemed that even slight changes to the network were changing the behavior greatly.

Further finetuning the hyperparams may provide some even more refined results. Changing the sigma coefficient to 0.02 helped the neural network at later episodes, and maybe an even lower value could help it even further. In possible future implementation I would decay the noise, starting from a value of around 0.05 all the way down to 0.01.

More algorithms for the same environment, like PPO, A3C, D4PG could possibly converge faster.