# Udacity Deep Reinforcement Learning Nanodegree

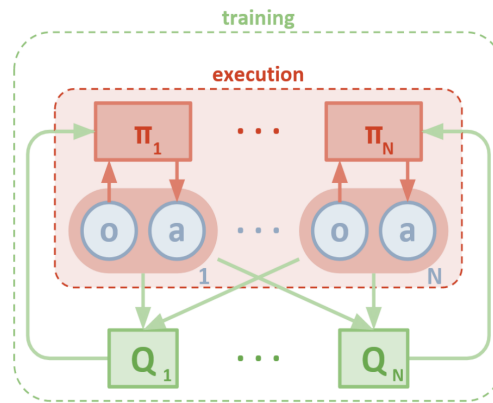## Project #3 Collaboration and Competition

## Report

### 1. Learning Algorithm

The algorithm implemented is the *MADDPG (Multiple Agent Deep Deterministic Policy Gradient) algorithm*, which can be considered as an extension of the DDPG algorithm, which in turn is a particular type of Actor-Critic method.

In the DDPG algorithm, in fact, we have one Critic part, which is used to approximate the maximizer over the Q values of the next states, which, in the case of continuous action spaces, would be too expensive to be calculated at each step of the agent. In addition to this, the MADDPG enables the so called decentralized-learning centralized-execution framework, where the Critic is trained with the observations, actions and possibly other information from all the agents of the environment. This, in addition to the fact that each agent is training its own Critic, means that is possible, for each agent, in both collaborative and competitive situations, to learn the best behavior taking into consideration the presence and behavior of other agents aswell. Once trained, the MADDPG needs to look only at the state of the environment, hence the centralized-execution part.

Each agent of the environment has 2 neural networks, one for the Actor, one for the Critic:

- The Actor approximates the optimal policy deterministically (i.e. it outputs the best believed action for any given state) and is it basically trying to learn the argmax over all actions for any state.

- The output action of the Actor, in conjunction with the state, aswell as the states and actions from all the other agents, are used to calculate a new target value for training the action-value function of the Critic of the agent.

  - The Critic, hence, is learning to evaluate the optimal value function by using the Actor best believed actions, always taking into considerations other agents actions, observations and possibly other information.

*Overview of the decentralized actor, centralized critic approach as explained in the
original MADDPG paper.*

## 2. Implementation

In this implementation, the MADDPG trains 2 agents, each with 2 neural networks:

- Actor Local (4 layers):
    - 1 with size equal to the state space (24 variables for the test environment)
    - 1 fully connected layer with 256 units
    - 1 fully connected layer with 128 units
    - 1 final layer with size equal to the action space (2 in the test environment)
        - The output of the final layer is then transformed with the Tanh function in order to get the actual action
- Actor Target
    - Same as Actor Local
- Critic Local (4 layers):
    - 1 with size equals to the state space for all the agents (24 * 2 agents in the test environment)
    - 1 fully connected layer with 256 units + the size of the action space of each agent (2 * 2 agents in the test environment)
    - 1 fully connected layer with 128 units
    - 1 final layer with size equal to one (the believed Q-value)
- Critic Target
    - Same as Critic Local

The Target versions of the neural networks are necessary to implement the soft update strategy, where, instead of updating the neural network after a fixed number of time steps, they are instead updated every time step by blending in the weights into the target network with a smoothing parameter governed by an hyperparameter, which has been shown to improve the stability of the learning process.

In order to guarantee the exploration of the environment, a noise factor, using the Ornstein-Uhlenbeck process, is applied to the actions returned by the Actor network.

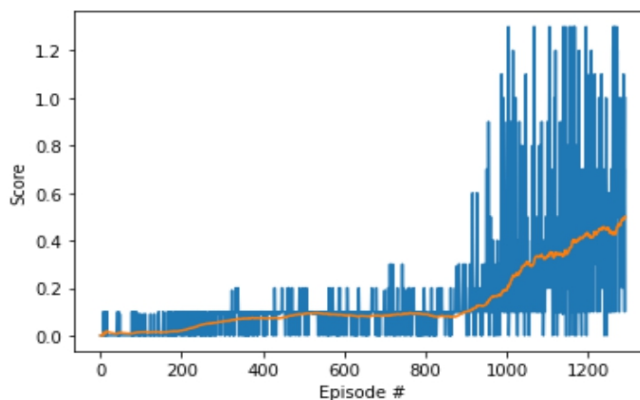As activation function between the layers, the ReLu function has been chosen.

The use of a Replay Buffer is necessary to learn from samples that are independent and identically distributed. Also, the learning does not happen at every time step but, instead, it happens after a fixed number of time steps and more than one update pass is applied at once (all regulated by hyperparameters).

## 3. Hyperparameters

For solving the environment, the hyperparameters have been chosen as following:

```python
BUFFER_SIZE = int(1e5)   # replay buffer size
BATCH_SIZE = 256         # minibatch size
GAMMA = 0.99             # discount factor
TAU = 0.09               # for soft update of target parameters
LR_ACTOR = 0.001         # learning rate of the actor
LR_CRITIC = 0.001        # learning rate of the critic
WEIGHT_DECAY = 0         # L2 weight decay
LEARN_EVERY = 4          # Learn every n steps
LEARN_HOWMANY = 5        # Learn how much every time
```

## 4. Plot of Rewards



*Plot of rewards over episodes.*

The environment has been solved (i.e. get *an average reward of the maximum between the the agents of +0.5 over 100 consecutive episodes*) in **1291** episodes. In blue the score for each episode (i.e. the maximum reward between the agents), in orange the moving average (i.e. the average of scores over the 100 previous episodes)

## 5. Ideas for Future Work

There are some possible improvements for the MADDPG algorithm:

1. Prioritized Experience Replay: sample experience transitions *uniformly* from a replay memory. [Prioritized experienced replay](#) is based on the idea that the agent can learn more effectively from some transitions than from others, and the more important transitions should be sampled with higher probability.

2. Other environments: It would be nice to experiment on other environments with the same MADDPG implementation in order to be able to compare the results, such as the proposed Soccer environment challenge.

3. The use of pixels as input would be a nice to have feature, updating the Neural Network in order to accommodate some initial Convolutional Layers as needed.

4. The use of parallel environments to speed up the training part.

## 6. Sources

- [https://github.com/udacity/deep-reinforcement-learning/tree/master/p3_collab-compet](https://github.com/udacity/deep-reinforcement-learning/tree/master/p3_collab-compet)

- [https://deepai.org/machine-learning-glossary-and-terms/relu](https://deepai.org/machine-learning-glossary-and-terms/relu)

- [https://www.udacity.com/course/deep-reinforcement-learning-nanodegree--nd893](https://www.udacity.com/course/deep-reinforcement-learning-nanodegree--nd893)

- [https://papers.nips.cc/paper/2017/file/68a9750337a418a86fe06c1991a1d64c-Paper.pdf](https://papers.nips.cc/paper/2017/file/68a9750337a418a86fe06c1991a1d64c-Paper.pdf)