

## Contents

DDPG Model implementation.....	2
Description of (MA)DDGP model use for collaboration game.....	2
Strategy used .....	2
Pseudo Code of the DDPG Algorithm: .....	2
Agent Parameters: .....	3
Parameters update .....	3
Agent's Actor and Critic network:.....	4
The Actor's parameters: .....	4
The Critic's parameters: .....	4
Dropout Rate.....	4
Batch normalization of actor and critic:.....	4
Random Noise generator: .....	4
Results of Model Training: .....	5
Graph of results .....	5
Potential Improvements: .....	5

## DDPG Model implementation

### Description of (MA)DDGP model use for collaboration game

The goal of this (MA)DDGP model is to train two Agents who bounce a tennis ball together. The game is collaborative. The (MA)DDPG is in parentheses because it is a wrapper around 2 agents so they share a common memory.

### Strategy used

The plan is to create 2 DDPG agents with common shared memory. Each agent is composed of a first DNN (the Actor) whose task is to approximate the action value function  $Q(s,a)$  by minimizing the squared distance of the Actor's predicted action/value vs its equivalent Bellman Equation while also training a second DNN (the Critic) whose task is to evaluate a given situation  $Q(s)$ .

The DDPG parameters of each agent had to be modified from assignment 2. The numbers of parameters in both actors & critic were increased to 512-256 for each of the first and second layer. A dropout rate was also added equivalent to 20%. Other dropout rates and smaller number of nodes were tried but were not as effective.

Also, the DDPG function was modified to accommodate the fact the agent receives multiple states at once. Similar to assignment two when 20 brains could be used to train for the task.

The unity environment provides 2 states at every step, one for each agent. The each agent is assigned an action corresponding to one of the 2 states. He then loops through the actions to store them in memory for future training steps into the common memory of the MADDPG which acts as a wrapper.

### Pseudo Code of the DDPG Algorithm:

1. Run a state input through each actor agent A and B
2. Collect experiences consisting of (actions, rewards, inputs) tuples from the actor
3. When enough experiences are memorized the update Actor & Critic Agent
4. Update critic agent
  - a. Use actor predicted next action to obtain  $Q_{next}$  as  $Critic(action_{next})$ 
    - i. Calculate target  $Q(S,A)$  as  $reward + discount * Q_{next}$
  - b. Calculate the expected  $Q(S,A)$  as  $Critic(action)$
  - c. Calculate Cost function as "target – expected"
5. Update the Agent's parameters:
  - a. For a given state sampled, predict the actions using the Actor
    - i.  $Action_{predicted} = Actor(state)$
  - b. Loss is defined as the average expected Critic Agent value
    - i.  $Average\_over\_actions(Critic(state, action_{predicted}))$
6. Perform as soft update of the Actor and critic target
  - a. Target actor & critic are update as  $x * target + (1-x) * local\_model$

### Agent Parameters:

***BUFFER\_SIZE = int(1e6) # replay buffer size***

***BATCH\_SIZE = 512***    *# minibatch size*

***GAMMA = 0.99***      *# 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*

## Parameters update

The parameters of the Agent are updated 10 times every 5 steps. The code that performs this is as follow. I had to extract some of the code of the Agent and put it into the DDPG function because otherwise the time step update was out of sync because of the multiple actions that are stored were previously counted as individual time step in my previous assignment. The “learnNiceThings() function is inside the MADDPG Agent ] This breaks down the step function from my previous assignment into two functions which is better because in OOP each function should have only one job. Only the learn Function is called within the singular DDGP agents

***# Learn, if enough samples are available in memory***

```
def learNiceThings(self):
```

```
self.t_step = (self.t_step + 1) % UPDATE_EVERY
```

```
if self.t_step == 0:
```

***# If enough samples are available in memory, get random subset and learn***

```
if len(self.memory) > BATCH_SIZE:
```

***for agent in self.ddpg\_agents:***

```
for _ in range(NUM_UPDATES):
```

```
experiences = self.memory.sample()
```

```
agent.learn(experiences, GAMMA)
```

## Agent's Actor and Critic network:

The model's architecture for approximating the  $Q(s,a)$  is the one described by the Google Deep Mind research paper.

### The Actor's parameters:

The actor consist two hidden layer fully connected neural network of with one hidden layer of 512-256 nodes and an output layer of 2 nodes describing the action value function for each 2 possible actions. The output uses the tanH so that the model's output are between -1 and 1

### The Critic's parameters:

The critic consist of a DNN fully connected with 2 hidden layers of 512-256 nodes reaching a final output of 1 that is between 0 and 1 to value the current state being passed to him.

The critic model also has a gradient clipping to improve performance

```
torch.nn.utils.clip_grad_norm(self.critic_local.parameters(), 1)
```

### Dropout Rate

A 20% dropout rate has been added to avoid overfitting. It has been tried in other cases to uses dropouts, it almost always have a positive effects. 20% is a standard starting point.

### Batch normalization of actor and critic:

The input of both the actor and the critic use batch normalization.

This has for effect to dampen the volatility that would be cause by an input with high low values. Example one dimension of the input is between -1000,1000 and the other one between 0-1, the first one will dominate the model and training will be poor.

### Random Noise generator:

The Theta parameter has been changed from 0.2 to 0.05. Also, The original code generated numbers between 0-1 but has been changed to the following:

```
#dx = self.theta * (self.mu - x) + self.sigma * np.array([random.random() for i in range(len(x))])
```

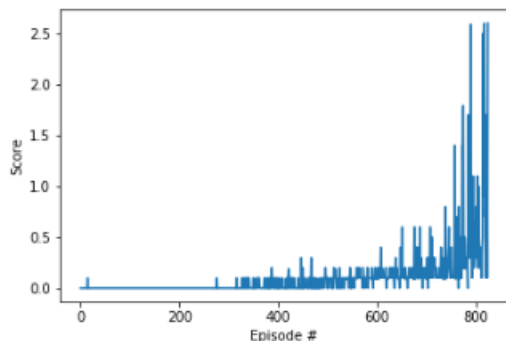
```
dx = self.theta * (self.mu - x)+self.sigma*np.random.uniform(-1, 1, len(x))
```

## Results of Model Training:

The problem is solved with roughly 825 episodes with a (MA)DDPG with augmented number of nodes and dropout rates. This again surprises me given that with slightly different parameters it take at least twice as much time or maybe even not at all.

### Graph of results

Episode 25	Average Score: 0.00
Episode 50	Average Score: 0.00
Episode 75	Average Score: 0.00
Episode 100	Average Score: 0.00
Episode 125	Average Score: 0.00
Episode 150	Average Score: 0.00
Episode 175	Average Score: 0.00
Episode 200	Average Score: 0.00
Episode 225	Average Score: 0.00
Episode 250	Average Score: 0.00
Episode 275	Average Score: 0.00
Episode 300	Average Score: 0.00
Episode 325	Average Score: 0.00
Episode 350	Average Score: 0.01
Episode 375	Average Score: 0.02
Episode 400	Average Score: 0.03
Episode 425	Average Score: 0.04
Episode 450	Average Score: 0.04
Episode 475	Average Score: 0.04
Episode 500	Average Score: 0.05
Episode 525	Average Score: 0.06
Episode 550	Average Score: 0.06
Episode 575	Average Score: 0.08
Episode 600	Average Score: 0.09
Episode 625	Average Score: 0.10
Episode 650	Average Score: 0.12
Episode 675	Average Score: 0.12
Episode 700	Average Score: 0.15
Episode 725	Average Score: 0.17
Episode 750	Average Score: 0.19
Episode 775	Average Score: 0.24
Episode 800	Average Score: 0.34
Episode 825	Average Score: 0.48
Episode 826	Average Score: 0.51\ Merry Christmas!! nEnvironment solved in 821 episodes! Average Score: 0.51



### Potential Improvements:

Using MADDPG could improve overall performance but it is a more complex algorithm to calibrate. I tried but it didn't work really well so that is why I reverted to simpler.

The hardest part of developing DRL algorithms is that you don't know in advance what will work in terms of number of layers, calibration parameters and all other fine tuning. The only way to know is to let it run for a long time and see what happens. Maybe the best way to circumvent this problem is people sharing their work so that this "evolution trial and error" happens as fast as possible and that the same mistakes are not repeated too often by too many people. As an example, the approach I took to use 2 agents is weaker than using only 1 agent when using 512-256 nodes and I don't know why but it is.

It certainly looks to me as more art than science for now.

Also, there is certainly a better way to code this whole thing than what I have done. Although it works, it doesn't appear very scalable because of the hard coded parameters.