**Project: Collaboration and Competition**

**Name: Zain Us Sami Ahmed Ansari**

## Introduction

For this project, I trained an agent to work with the Tennis environment.

In this environment, two agents control rackets to bounce a ball over a net. If an agent hits the ball over the net, it receives a reward of +0.1. If an agent lets a ball hit the ground or hits the ball out of bounds, it receives a reward of -0.01. Thus, the goal of each agent is to keep the ball in play.

The observation space consists of 8 variables corresponding to the position and velocity of the ball and racket. Each agent receives its own, local observation. Two continuous actions are available, corresponding to movement toward (or away from) the net, and jumping.

## Implementation

In this exercise the parameters that are tuned are as follows.

```
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 = 5e-4          # learning rate of the actor
LR_CRITIC = 1e-4         # learning rate of the critic
WEIGHT_DECAY = 0.0000    # L2 weight decay
```

The model architecture is defined through set of following variables

```
Number of agents: 2
Size of each action: 2
There are 2 agents. Each observes a state with length: 24
```

The parameters that define the network are as follows

**"""Actor (Policy) Model."""**

> """Initialize parameters and build model.
> Params
> ======
> > state_size (24): Dimension of each state
> > action_size (2): Dimension of each action
> > seed (1): Random seed
> > fc1_units (128): Number of nodes in first hidden layer
> > fc2_units (64): Number of nodes in second hidden layer
> """

**"""Critic (Value) Model."""**

    """Initialize parameters and build model.
    Params
    ======
        state_size (24): Dimension of each state
        action_size (2): Dimension of each action
        seed (1): Random seed
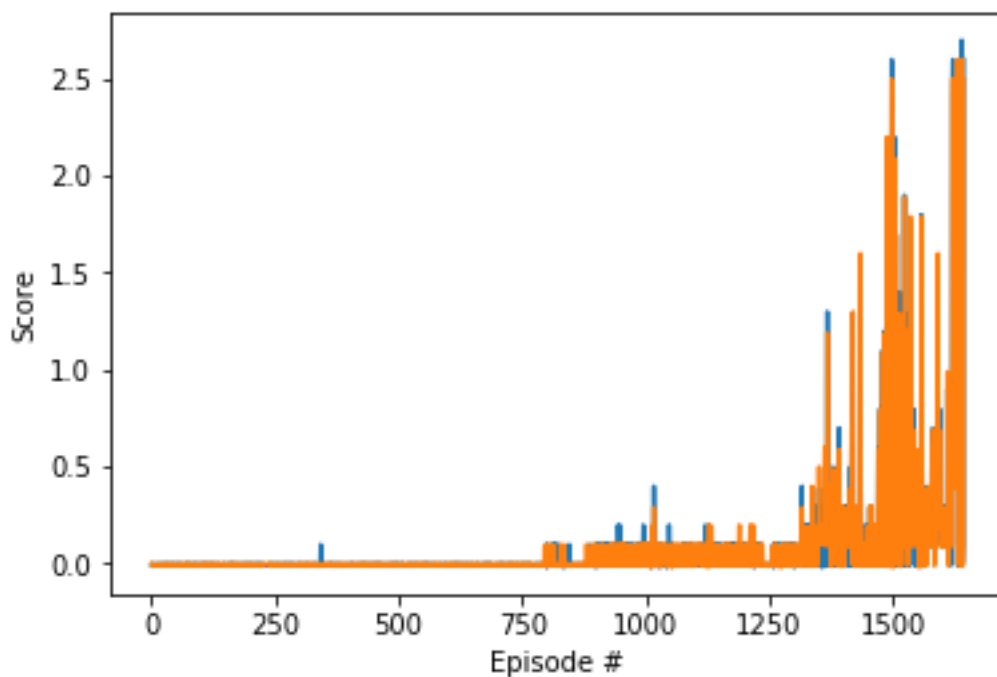        fc1_units (128): Number of nodes in the first hidden layer
        fc2_units (64): Number of nodes in the second hidden layer
    """

I have used two neural networks, the actor network and the critic network. The actor network contain 2 hidden layers of 128 and 64 units respectively.  The output layer is a tanh activation layer.

The critic network has two hidden layers 128 and 64 units respectively.

## Plot of Rewards



```
Environment solved in 1544 episodes!   Average Score:
0.509800
```

## Future Work

As it is very evident from the plot of rewards that the standard deviation of this implementation is very high I want to try out different parameters and network architectures to reduce the standard deviation.

I would also want to try out algorithms like PPO, A3C and D4PG that use multiple (non-interacting, parallel) copies of the same agent to distribute the task of gathering experience.
I would also like to try changing different actor, critic network sizes in DDPG.