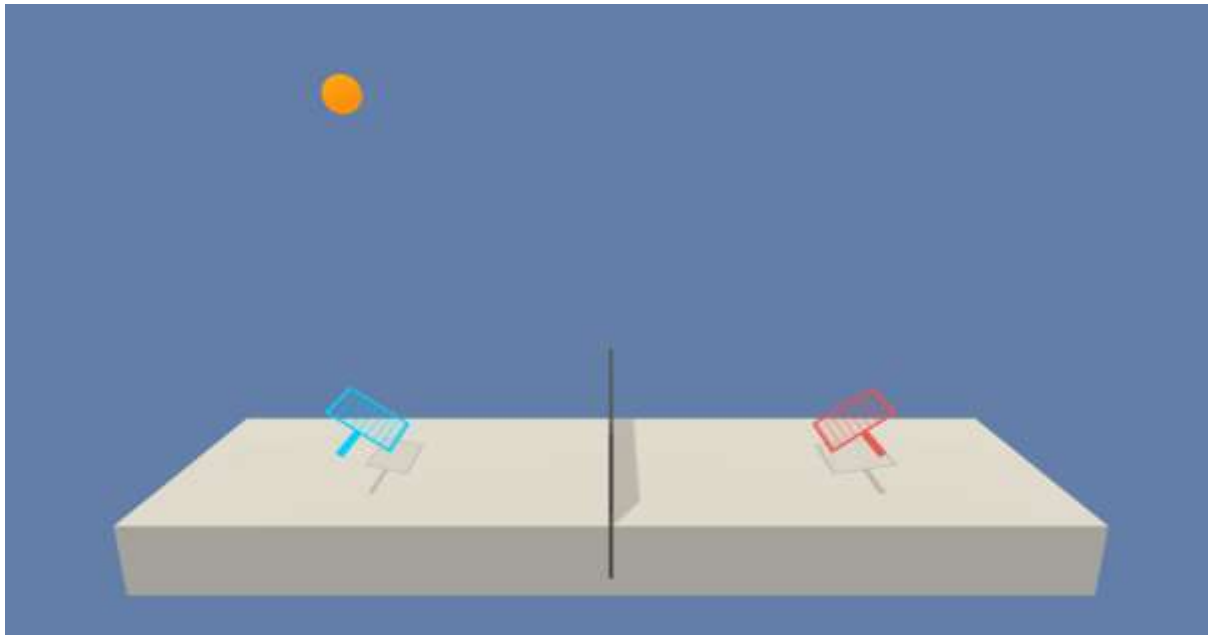


Project 3 : Collaboration-Competitions

For this project, you will work with the **Tennis** environment.



Unity ML-Agents 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.

The task is episodic, and in order to solve the environment, your agents must get an average score of +0.5 (over 100 consecutive episodes, after taking the maximum over both agents). Specifically,

- After each episode, we add up the rewards that each agent received (without discounting), to get a score for each agent. This yields 2 (potentially different) scores. We then take the maximum of these 2 scores.
- This yields a single **score** for each episode.

The environment is considered solved, when the average (over 100 episodes) of those **scores** is at least +0.5.

Algorithm

Learning Algorithm:

Multi-Agent Deep Deterministic Policy Gradient Algorithm

For completeness, we provide the MADDPG algorithm below.

Algorithm 1: Multi-Agent Deep Deterministic Policy Gradient for N agents

```
for episode = 1 to  $M$  do
  Initialize a random process  $\mathcal{N}$  for action exploration
  Receive initial state  $\mathbf{x}$ 
  for  $t = 1$  to max-episode-length do
    for each agent  $i$ , select action  $a_i = \mu_{\theta_i}(o_i) + \mathcal{N}_t$  w.r.t. the current policy and exploration
    Execute actions  $\mathbf{a} = (a_1, \dots, a_N)$  and observe reward  $r$  and new state  $\mathbf{x}'$ 
    Store  $(\mathbf{x}, \mathbf{a}, r, \mathbf{x}')$  in replay buffer  $\mathcal{D}$ 
     $\mathbf{x} \leftarrow \mathbf{x}'$ 
    for agent  $i = 1$  to  $N$  do
      Sample a random minibatch of  $S$  samples  $(\mathbf{x}^j, \mathbf{a}^j, r^j, \mathbf{x}^{\prime j})$  from  $\mathcal{D}$ 
      Set  $y^j = r^j + \gamma Q_{\theta_i'}^{\mu}(\mathbf{x}^j, \mathbf{a}^j, \dots, \mathbf{a}_N^j) \big|_{\mathbf{a}_i' = \mu_{\theta_i'}^{\mu}(\mathbf{o}_i^j)}$ 
      Update critic by minimizing the loss  $\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j (y^j - Q_{\theta_i}^{\mu}(\mathbf{x}^j, \mathbf{a}_1^j, \dots, \mathbf{a}_N^j))^2$ 
      Update actor using the sampled policy gradient:
      
$$\nabla_{\theta_i} J \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(\mathbf{o}_i^j) \nabla_{\mathbf{a}_i} Q_{\theta_i}^{\mu}(\mathbf{x}^j, \mathbf{a}_1^j, \dots, \mathbf{a}_i, \dots, \mathbf{a}_N^j) \big|_{\mathbf{a}_i = \mu_{\theta_i}(\mathbf{o}_i^j)}$$

    end for
    Update target network parameters for each agent  $i$ :
    
$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$

  end for
end for
```

It uses two Target networks to add stability to training

To implement better exploration by the Actor network, we use noisy perturbations, specifically an Ornstein-Uhlenbeck process for generating noise. It uses samples noise from a correlated normal distribution.

Critic loss - Mean Squared Error of $y - Q(s, a)$ where y is the expected return as seen by the Target network, and $Q(s, a)$ is action value predicted by the Critic network. y is a moving target that the critic model tries to achieve; we make this target stable by updating the Target model slowly.

Actor loss - This is computed using the mean of the value given by the Critic network for the actions taken by the Actor network. We seek to maximize this quantity.

Neural network architecture:

state_size: Dimension of each state

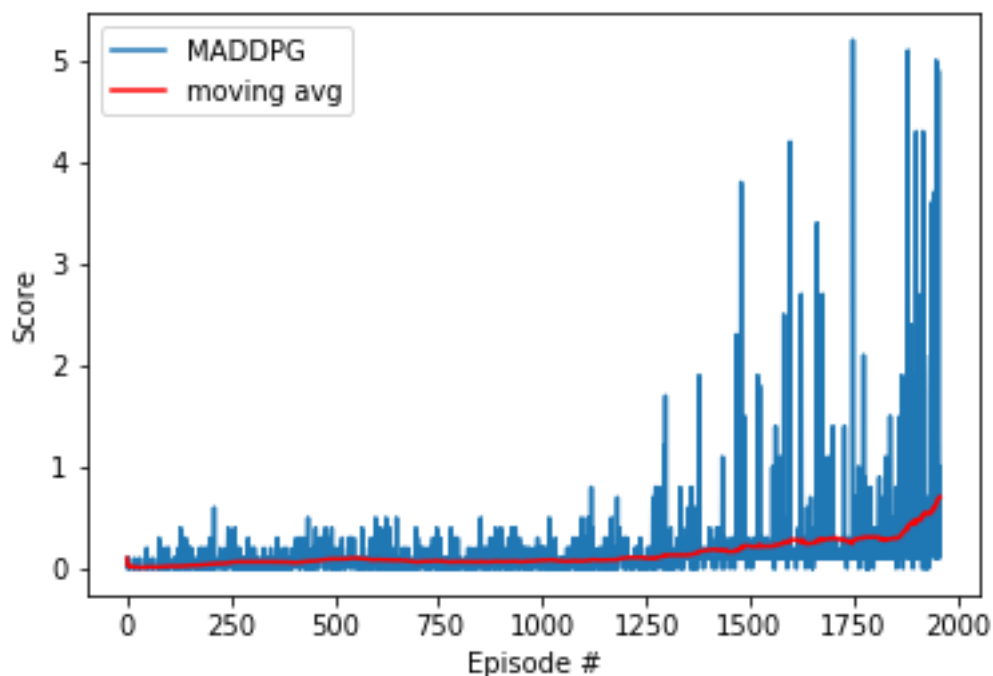
action_size: Dimension of each action

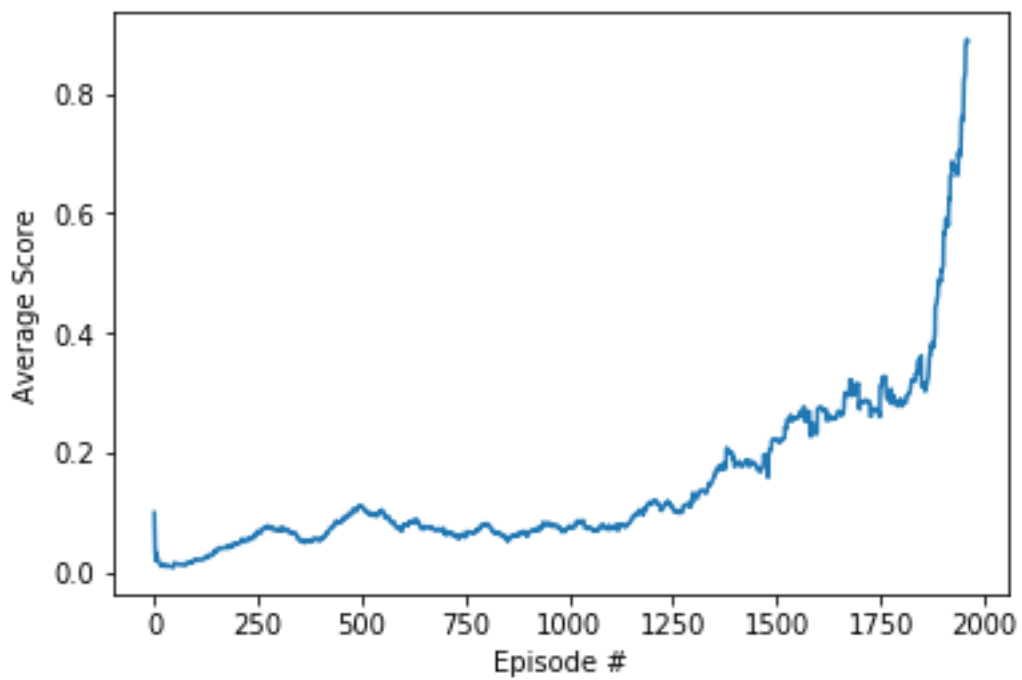
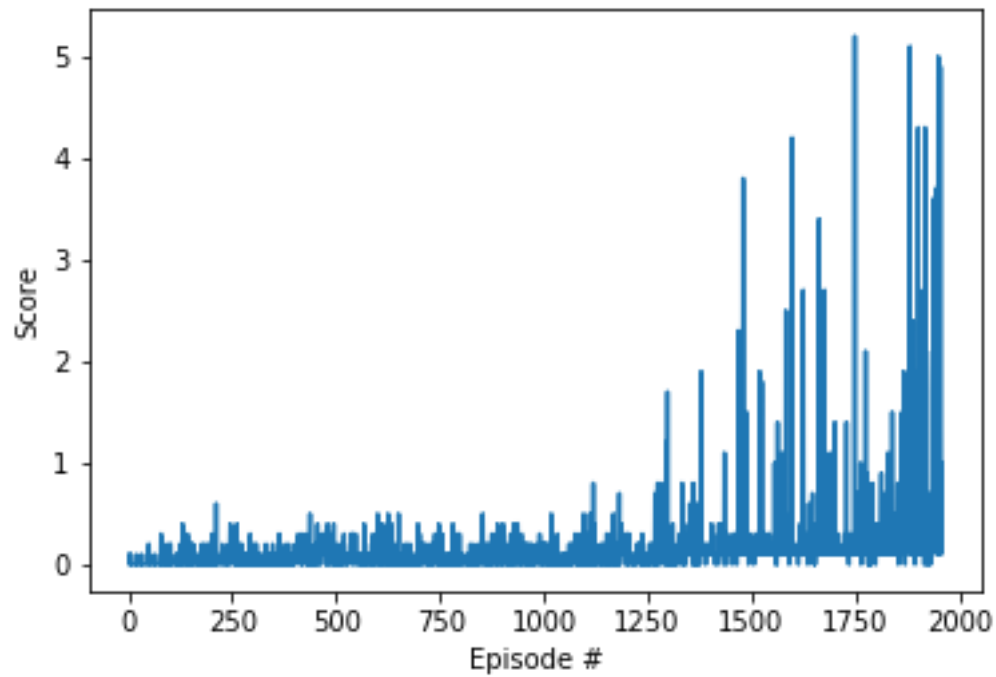
fc1_units: int. Number of nodes in first hidden layer – 128 Units
fc2_units: int. Number of nodes in second hidden layer – 128 Units
We set same for both actor and critic network

Other parameters are:

```
BUFFER_SIZE = int(1e6) # replay buffer size
BATCH_SIZE = 128      # minibatch size
LR_ACTOR = 1e-3        # learning rate of the actor
LR_CRITIC = 1e-3       # learning rate of the critic
WEIGHT_DECAY = 0       # L2 weight decay
LEARN_EVERY = 1        # learning timestep interval
LEARN_NUM = 5          # number of learning passes
GAMMA = 0.99           # discount factor
TAU = 8e-3             # for soft update of target parameters
OU_SIGMA = 0.2         # Ornstein-Uhlenbeck noise parameter, volatility
OU_THETA = 0.15        # Ornstein-Uhlenbeck noise parameter, speed of mean reversion
EPS_START = 5.0        # initial value for epsilon in noise decay process in Agent.act()
EPS_EP_END = 300       # episode to end the noise decay process
EPS_FINAL = 0          # final value for epsilon after decay
```

Here is the training graph





The environment is considered solved, when the average (1890-1900episodes) of those average scores is at least 0.5.

Future with multiagent:

GAE - Generalized Advantage Estimation

A3C - Asynchronous Advantage Actor-Critic

A2C - Advantage Actor-Critic

PPO - Proximal Policy Optimization

D4PG - Distributed Distributional Deterministic Policy Gradients

References :

<https://keras.io/examples/rl/>