

Description :

The challenge requires to train 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. 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).

Learning Algorithm : MADDPG

The learning algorithm is Multi Agent Deep Deterministic Policy Gradient which is an improvement of the solution used in Bipedal Gym environment in the class. The agent is unaware of the entire state and tries to get reward through following policy based algorithm. The architecture employed to solve this environment is as follows:

Actor Network:

3 Fully connected layers with ReLU activations and tanh on the last layer. Learning Rate: $1e-4$

The network is normalized through dropout and batch normalization.

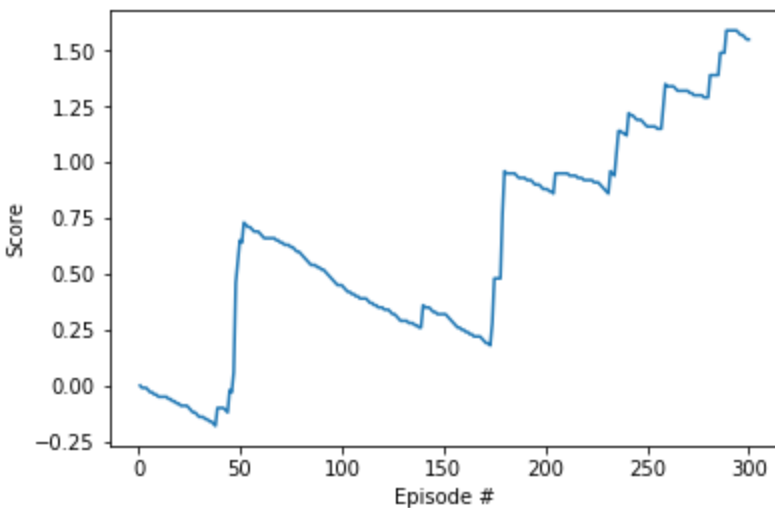
Critic Network:

3 Fully connected layers with ReLU activations. Learning Rate: $1e-3$

The network is normalized through dropout.

Parameters: Replay buffer size : 1000000, Minibatch size: 1024, $\tau = 1e-3$, Discount factor: 0.99

Agent and Environment: The environment is solved in 210 episodes.



The algorithm seems hard to converge without normalization. With normalizations such as batch normalization and dropout applied to the networks, the algorithm converges faster.

Future Work:

DDPG algorithm works quite well in this environment. Further, hyperparameter tuning would have given improvements given the tennis environment is solved in few hundred episodes, we can confidently try PPO, TRPO and actor critic methods such as A3Cs. I will follow these github repositories to improve upon my knowledge and modularity

<https://github.com/qfettes/DeepRL-Tutorials>

<https://github.com/ShangtongZhang>

As always I will watch out for recent trends from OPENAI blog. :)
