

Udacity Deep Reinforcement Learning

Project: Multi Agent Reinforcement Learning

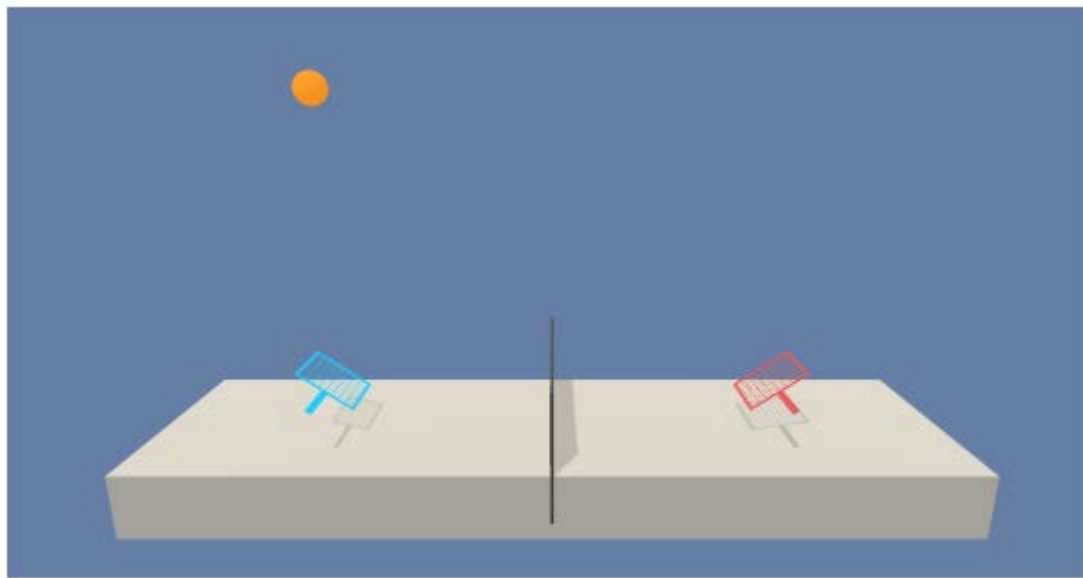
Author – Anish Amul Vaidya (avaidya172@gmail.com)

Abstract - 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.

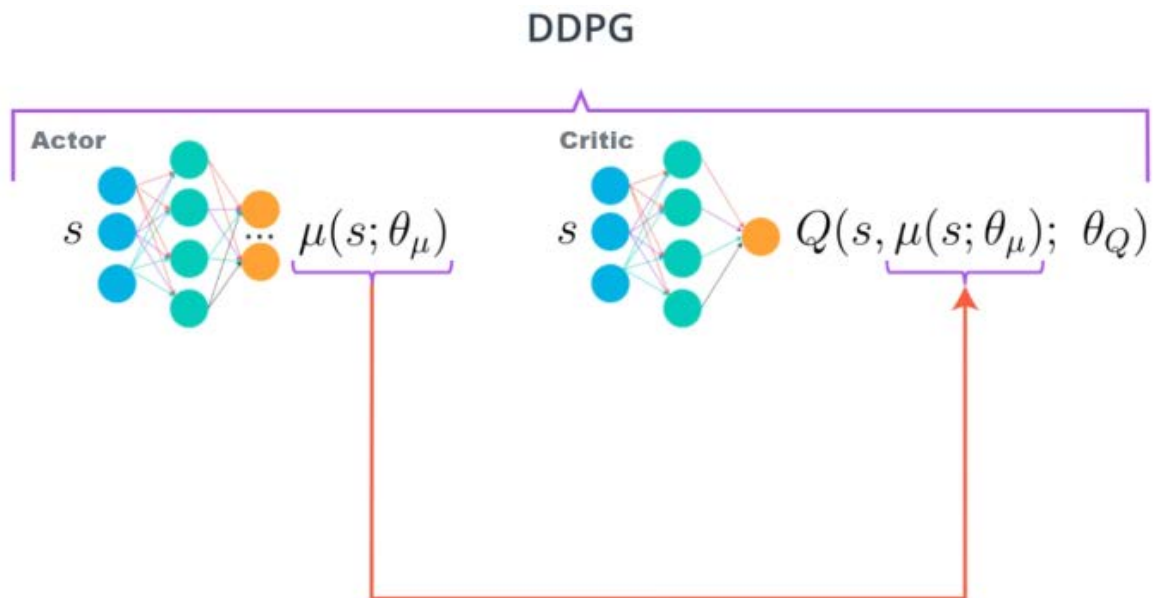


The environment is considered solved, when the average (over 100 episodes) of those scores is at least +0.5. The following report is written in four parts:

- Implementation
- Results
- Possible future improvements

Implementation

DDPG(Lilicrapetal.,2015), short for Deep Deterministic Policy Gradient, is a model free off-policy actor-critic algorithm, combining DPG with DQN. Recall that DQN (Deep QNetwork) stabilizes the learning of the Q-function by using experience replay and a frozen target network. The original DQN works in discrete space, and DDPG extends it to continuous space with the actor-critic framework while learning a deterministic policy. In DDPG, we use 2 deep neural networks : one is the actor and the other is the critic:

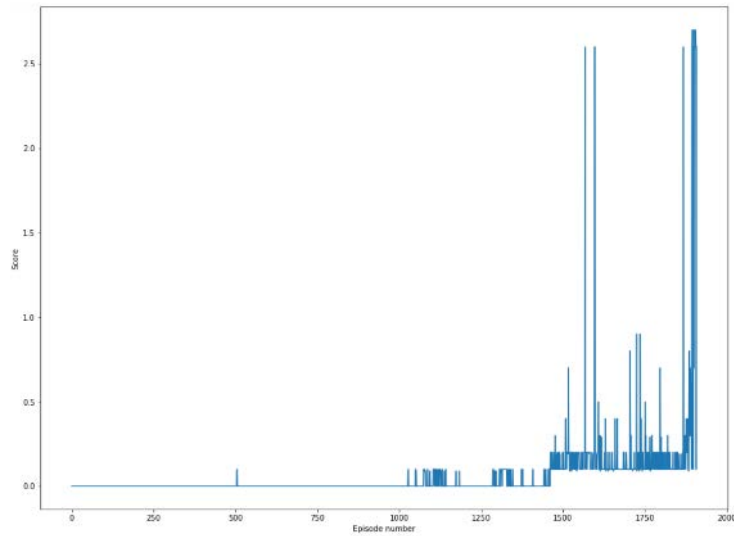


There is a modification to the DDPG agent used in the 'continuous control - robot arm' project. The main difference is that here, all the agents share the same replay buffer memory. The main reason that the agent may get caught in a tight loop of specific states that loop back through one another and detract from the agent exploring the full environment "equally". By sharing a common replay buffer, we ensure that we explore the whole environment. Both the agents still have their own actor critic networks to train.

Because we have two models for each agent; the actor and critic that must be trained, it means that we have two set of weights that must be optimized separately. Adam was used for the neural networks with a learning rate of 10^{-4} and 10^{-3} respectively for the actor and critic for each agent. * For Q, I used a discount factor of $\gamma = 0.99$. For the soft target updates, the hyperparameter τ was set to 0.001. The neural networks have 2 hidden layers with 250 and 100 units respectively. For the critic Q, the actions were not included the 1st hidden layer of Q. The final layer weights and biases of both the actor and critic were initialized from a uniform distribution $[-3 \times 10^{-3}, 3 \times 10^{-3}]$ and $[3 \times 10^{-4}, 3 \times 10^{-4}]$ to ensure that the initial outputs for the policy and value estimates were near zero. As for the layers, they were initialized from uniform distribution $[-1/\sqrt{f}, 1/\sqrt{f}]$ where f is the fan-in of the layer. We use the prelu1 activation function for each layer.

Results

The following is the graph of the scores per episode as the models train:



Possible future improvements

Because of the symmetry of the environment, each agent mirrors the other one behaviour. As such, there is no really such adversarial or collaborative relationship needed, thus the multiagent reinforcement learning does not seem ideal in this setting. A great challenge would be to tackle an environment in which the Multi-Agent Reinforcement Learning framework could apply, like the Soccer environment. With it, we can apply Multi-Agent DDPG (Lowe et al.) in which each of the actor takes a concatenation of all the states and actions of the other agents.