

DDPG Algorithm - Reacher Continuous Control

Overview

Using the Unity agent/environment "Reacher", this deep reinforcement learning task trains an AI agent to keep its double jointed arm on a ball for as long as possible. The task is considered solved when the agent, who receives various rewards for maintaining its position, receives an average reward of 30 over 100 consecutive episodes.

The Unity environment has two options - using a single agent or 20 agents in parallel. In the latter case, the experiences and rewards of all 20 agents running at the same time are combined during evaluation and training. This repo uses the 20 agent version.

In either environment, each agent receives feedback in the form of a reward (+0.1) or no reward after taking each action. The environment conveys to the agent the state of the agent in the form of 33 different values related to position, rotation, velocity and angular velocities of the arm. Each action is a vector with four numbers, corresponding to torque applicable to two joints.

At first the agent randomly takes actions and records the feedback, but eventually begins to take those experiences and learn from them using a Deep Deterministic Policy Gradient (DDPG) algorithm.

The attached code written in Python, using PyTorch and presented in a Jupyter Notebook, demonstrates how the agent learns and eventually achieves an average score of 30 over 100 consecutive episodes.

Model Architecture

The Udacity provided DDPG code in PyTorch was used and adapted for this 20 agent (version 2) environment.

I used two deep neural networks (actor-critic), each with two hidden layers of 256-128 nodes, with eLU activation functions on the hidden layers and tanh on the output layers.

The changes I made vs the original DDPG research paper -

<https://arxiv.org/pdf/1509.02971.pdf> - was reducing network complexity, after experimenting I found it converges faster this way (from 400-300 to 256-128).

Also, in the original source code from Udacity the learning rates were different for the two networks, I found it better to use same for both $1e-4$.

Hyperparameters

A learning rate of $1e-4$ on each DNN and batch size of 128 were used along with replay buffer size of $1e5$, gamma .99 and Tau of $1e-3$. There was no change to the default Ornstein-Uhlenbeck noise parameters (0.15 theta and 0.2 sigma.)

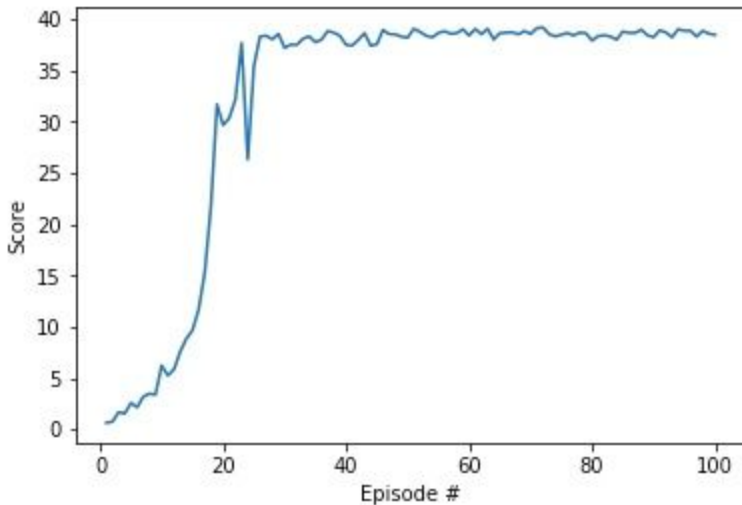
Training

I split the training into two parts: at first, since we are using 20 agents that share experiences, thus generating a lot of experiences, I train the agents 3 times per episode, until it reaches an acceptable value over the averages (33 chosen). After this level is achieved, I train each 10 episodes.

In addition to Ornstein-Uhlenbeck noise, I found helpful to add an epsilon greedy logic to the agent.

Results and Future Work

The model was able to achieve the 30 average reward goal around episode 20, and it maintained performance afterwards, thus achieving the goal in the first 100 episodes.



In the future, I plan experimenting with other architectures, such as TRPO, PPO, etc.