

Project criteria

The goal of this project was to train robotic arm to maintain contact with the sphere in the environment. The agents must get an average score of +30, over 100 consecutive episodes.

Implementation

The algorithm that I have used is an implementation of Deep Deterministic Policy Gradient with Replay Buffer. I have also found that features such as gradient clipping and learning every 10 episodes but 10 times helped to speed up training process and made it more stable.

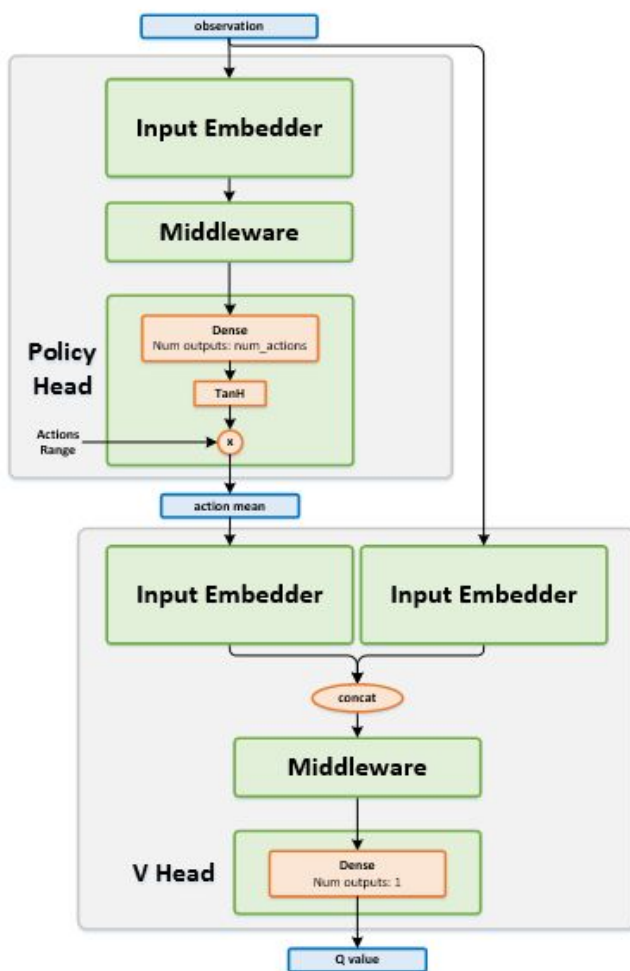
Original paper:

<https://arxiv.org/pdf/1509.02971.pdf>

Helpful explanation:

<https://towardsdatascience.com/deep-deterministic-policy-gradients-explained-2d94655a9b7b>

DDPG is bringing the best from worlds of Policy methods and Value methods by introducing Actor and Critic networks focused on different sides of the problem. Actor learns the action based on observation, while Critic network learns the Q value of the state and the action produced by Actor. With the addition of target Actor and target Critic networks, which are slightly delayed copies of originals, we balance instabilities and allow step by step movement of the entire system towards the solution of the environment.



Actor network:

- 3 layer neural network with additional batch normalization layer
- 33 -> 400 -> Batch norm -> 300 -> 4

Critic network:

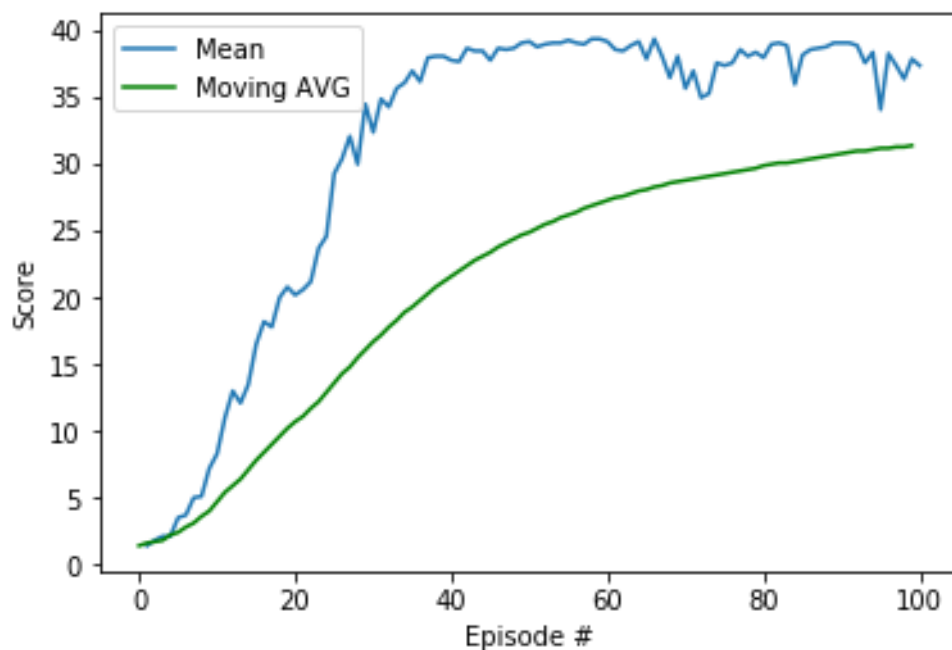
- 3 layer neural network with additional batch normalization layer
- 33 -> 400 -> Batch norm -> 300 -> 4

Hyperparameters:

- replay buffer size = 1e6
- minibatch size = 128
- discount factor = 0.99
- tau = 0.0025
- learning rate of the actor = 0.0002
- learning rate of the critic = 0.0005

- L2 weight decay = 0
- epsilon = 1
- epsilon decay = 1e-6
- epsilon min 0.1

Results



The agents were able to solve task in 100 episodes with a final moving average score of 31.4.

Improvements

- Experimenting with different RL algorithms may yield more stable training process (D3PG, D4PG, A3C, PPO) as this is one of the main flaws of DDPG.
- Improving results tuning the hyperparameters may decrease instability of the DDPG algorithm.
- Different Actor and Critic network structures could help with more efficient features extraction.
- Prioritized experience replay could be very useful to learn from situations with high magnitude of error.