

For this project, I implemented the DDPG Actor-Critic algorithm for continuous action spaces. It was the right algorithm to implement due to its combination of policy gradient and value function approximation. The combination of the two lead to better stability and results than either algorithm alone. Below are the steps for implementation.

1. Initialize replay memory buffer
2. Initialize local actor and critic network as pytorch neural nets
3. Implement a soft-weights copy function
4. Train the actor and the critic

The actor and critic network have 3 fully connected layers, with the only difference being a tanh function at the end of the actor network to transform the outputs between -1 and 1. One other minor but crucial difference is that the second layer of the critic network has an input size to include both the output of the first layer AND the actions chosen by the actor network. This enables the critic to evaluate the actor.

The first layer of the Actor network has an input size of 33 (observation size) by 264 (chosen number of output nodes): 33 x 264

The second layer of the Actor network has a size of 264 x 132(chosen number of output nodes for this layer): 264 x 132

The third layer of the Actor network has a size of 132 x 4 (action size): 132 x 4

The first layer of the Critic network has a size of 33(observation size) by 264 (chosen number of output nodes): 33 x 264

The second layer of the Critic network has a size of 268 (output nodes from previous layer + action size) by 132 (chosen number of output nodes for this layer): 268 x 132

The third layer has a size of 132 (output nodes from previous layer) by 1 (scalar for value of current state action pair): 132 x 1

The hyperparameters used for training the network are as follows:

`BUFFER_SIZE = int(1e5)`

`BATCH_SIZE = 1000`

GAMMA = 0.99

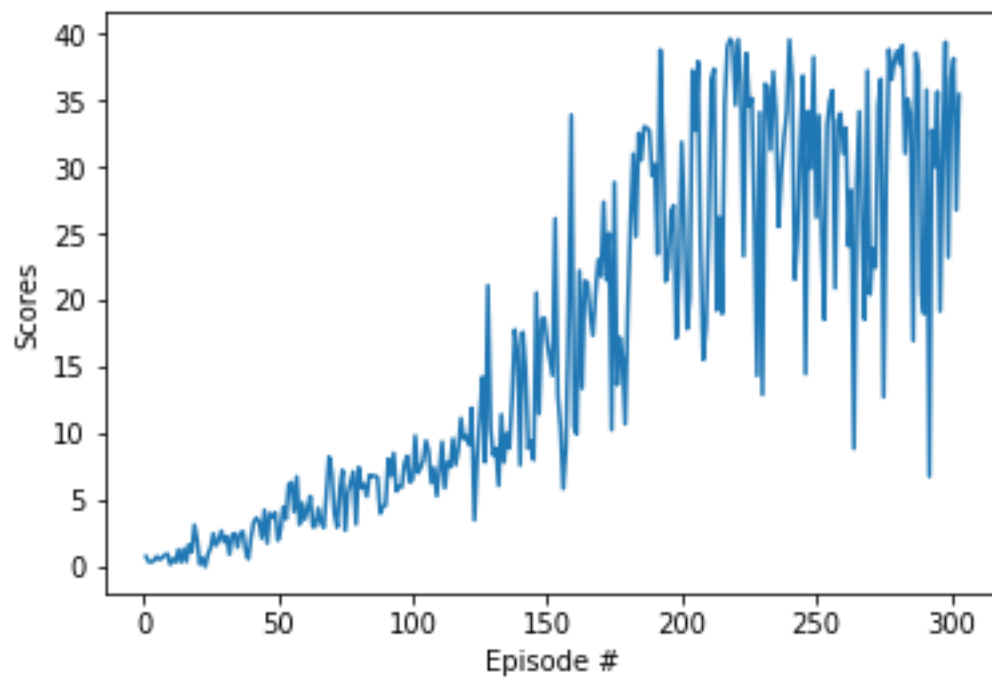
TAU = 1e-3

LR_ACTOR = 1e-4

LR_CRITIC = 1e-3

WEIGHT_DECAY = 0

Here is a plot of the rewards. As we can see, after the 200th episode the environment is considered solved.



Ideas for future work:

1. Prioritized experience replay
2. Asynchronous Actor Critic