

## Continuous Control Report

### Algorithm Description

The Deep Deterministic Policy Gradient (DDPG) method was used to solve this environment. This method uses a neural network to estimate the action value function,  $Q$ , and another neural network to estimate the policy. These are referred to as the critic and actor networks respectively.

The specific implementation I used has an epsilon and a sigma value to control exploration. Each time an action is chosen, a small noise value is added to it. This noise comes from a gaussian distribution with mean 0 and standard deviation of sigma. That noise value is then multiplied by epsilon. As training occurs, epsilon is decreased.

A replay buffer was also used to allow the agent the train on uncorrelated experiences each epoch.

A local and target network were used to stabilize training with a soft update TAU each time step. These networks are initialized with a hard update step so that the local and target networks start with the same random weights. The Temporal Difference (TD) method was used to calculate  $Q$  along with the target critic network. This approach increases bias but decreases noise allowing the network to train more reliably.

To further improve stability, the networks train 10 epochs every 10 time steps.

### Hyperparameters

The following hyperparameters were used in this solution:

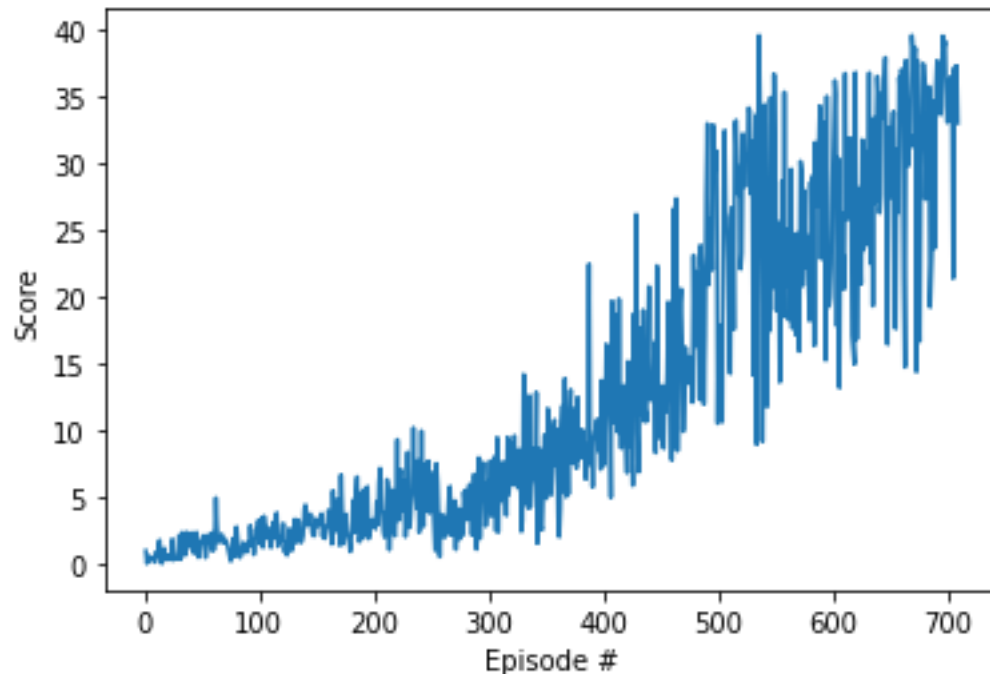
- BUFFER\_SIZE = int(1e6) – Size of the replay buffer
- BATCH\_SIZE = 128 – Number of experiences used in each training epoch
- GAMMA = 0.99 – Discount factor on rewards
- TAU = 1e-3 – Soft update interpolation parameter
- LR\_ACTOR = 1e-4 – Learning rate for actor network
- LR\_CRITIC = 3e-4 – Learning rate for critic network
- WEIGHT\_DECAY = 0.00001 – L2 weight decay on critic network
- eps\_start = 1.0 – Starting value of epsilon
- eps\_decay = 0.97 – Decay rate of epsilon ( $\text{eps}(t+1) = \text{eps}(t) * \text{eps\_decay}$ )

The actor network has 3 layers (1 input, 1 output, and 1 hidden). The input layer is the state size, the output layer is the action size, and the hidden layer has 64 nodes. Each layer uses a tanh activation function.

The critic network has 4 layers (1 input, 1 output, and 2 hidden). The input layer is the state size + action size, the output layer has 1 node (the predicted value), and the hidden layers all have 64 nodes. Each layer uses a tanh activation function except the output layer which has no activation function.

## Results

The agent took a total of 708 episodes to train. The score per episode can be seen in the plot below:



The final average reward over 100 episodes was 30.09.

## Future Work

This problem may be able to be solved more effectively using:

- Prioritized Experience Replay – Assign each experience an importance level based on its TD error then sample from the replay buffer non-uniformly where more important experiences are used more often.
- N step bootstrapping – Instead of using pure TD estimation or Monte Carlo estimation for the Q function, use a mixture of the two where n true rewards are used to compute the value. This increases noise but reduces bias.
- Modified Loss Function – The agent spent a long time during training spinning at high speeds. A modified loss function that penalizes high control input may help avoid that behavior and lead to more consistent performance during training.