

# Udacity DRLND Project #2 Continuous Control

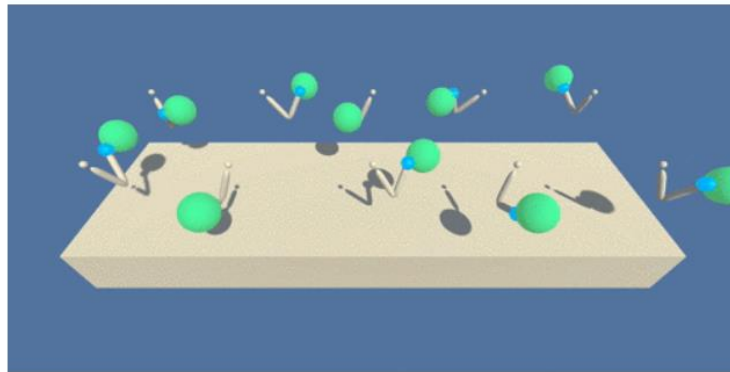
Jing Zhao 11/03/2018

## Introduction

In this report, I described in detail the implementation of a deep deterministic policy gradient (DDPG) algorithm for solving the continuous control (Reacher) problem in the Unity ML-Agents environment. I first introduced the environment followed by a deep dive into the learning algorithm. The environment is solved in 80 episodes. A plot of rewards per episode is showed to illustrate the training process, along with several key learnings to success which include large batch-size, infrequent update of network, and epsilon-decay of the noise injected into action space. A couple of ideas for future improvement are presented at the end.

## Environment

The Reacher environment is provided by DRLND's course website. In this environment, a double-jointed arm can move to target locations. A reward of +0.1 is provided for each step that the agent's hand is in the goal location. Thus, the goal of my agent is to maintain its position at the target location for as many time steps as possible. Figure 1 is a screenshot of the environment.



**Figure 1.** A screenshot of the Continuous Control (Reacher) environment.

The observation space consists of 33 variables corresponding 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. Every entry in the action vector should be a number between -1 and 1.

For this project, two separate versions of the Unity environment are provided:

- The first version contains a single agent.
- The second version contains 20 identical agents, each with its own copy of the environment.

There are two options to solve the environment:

*Option 1: Solve the First Version*

The task is episodic, and in order to solve the environment, my agent must get an average score of +30 over 100 consecutive episodes.

### Option 2: Solve the Second Version

The barrier for solving the second version of the environment is slightly different, to take into account the presence of many agents. In particular, my agents must get an average score of +30 (over 100 consecutive episodes, and over all agents).

### Learning Algorithm

I used the DDPG algorithm [1] to solve this environment. DDPG is basically an actor-critic method. The actor network maps states to the continuous action space, while the critic network approximates the Q-value of state-action pair. The goal is to maximize the Q-value by adjusting the weights of actor and critic networks. Inspired by the success of deep Q-network (DQN) [2], updating critic network is also driven by minimizing the temporal-difference from one time-step to the next. The two significant improvements of DQN over traditional online Q-learning - namely experience replay and fixed target network – are adopted in the DDPG algorithm. A pseudocode from [1] is shown below for reference.

---

#### Algorithm 1 DDPG algorithm

---

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .  
Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$   
Initialize replay buffer  $R$   
**for** episode = 1, M **do**  
  Initialize a random process  $\mathcal{N}$  for action exploration  
  Receive initial observation state  $s_1$   
  **for** t = 1, T **do**  
    Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise  
    Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$   
    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$   
    Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$   
    Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$   
    Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$   
    Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

Update the target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

**end for**  
**end for**

---

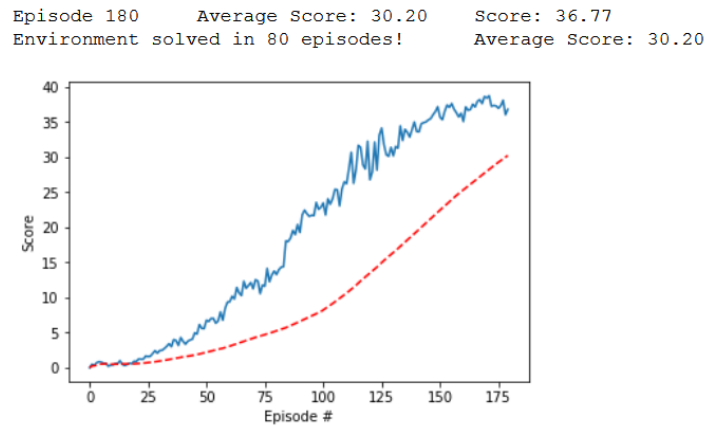
I built two 2-layer feedforward neural networks (NN) for the actor and critic, respectively. Table 1 shows the hyperparameters of these two networks along with some other learning parameters in the DDPG algorithm.

**Table 1.** Hyperparameters of DDPG and other learning parameters

Parameter	Description	Value
State_size	Dimension of the environment	33
Action_size	Dimension of the action space	4
BUFFER_SIZE	Replay buffer size	1,000,000
BATCH_SIZE	Minibatch size	1024
GAMMA	Discount factor	0.99
TAU	For soft update of target parameters	1e-3
LR_ACTOR	Learning rate of actor network	5e-4
LR_CRITIC	Learning rate of critic network	1e-3
UPDATE_EVERY	How often (timesteps) to update the target network	100
UPDATE_RUNS	If update, how many times in a row	10
N_episodes	Total episodes in training	2000
Max_t	Maximum time steps in each episode	1000
[state_size, fc1_units, fc2_units, action_size]	Units in each layer of the actor network	[33, 128, 64, 4]
[[state_size, action_size], fcs1_units, fc2_units, output_size]	Units in each layer of the critic network	[[33, 4], 64, 64, 1]
Eps_start	Start value of noise (exploration) decay	1.0
Eps_end	End value of noise (exploration) decay	0.1
Eps_decay	Rate of exponential noise (exploration) decay	0.995

## Result and Discussion

I leveraged the DDPG Python implementation from the course GitHub repository (ddpg-pendulum) for this project. Figure 2 below shows the score of each episode (in blue) as well as the moving-average across past 100 episodes (in red dash line). Note that these scores are averaged across 20 agents in the second version of environment. According to the project description, this environment is “solved” once the agent is able to receive an average reward (over 100 episodes, across all agents) of at least +30. I achieved this goal in 80 episodes.



**Figure 2.** Averaged score over 20 agents (blue) and its 100-point moving average (red) during training using a DDPG algorithm. The environment is solved in 80 episodes.

Next, I'd like to report a few crucial findings that are essential to my success.

1. **Batch Size.** By far, this is the parameter that dominates my agent's performance. I ended up using  $N=1024$  for the reason of performance stability and learning speed. The agent failed to learn anything if a much smaller batch (e.g.  $N = 64$ ) was used. With  $N=512$ , the environment was solved in 120 episodes with all other hyperparameters remained the same, indicating that doubling the batch size does accelerate learning.
2. **Frequency of Update.** Infrequent update of the target network stabilizes my agent's performance quite a bit. I eventually landed on a policy that updates my networks 10 times every 100 time-steps. It also brings in some efficiency in training.
3. **Network Size.** I used two 2-layer feedforward networks with each layer having no more than 128 units. I feel comfortable with my network size. From my readings that the reinforcement learning community also recognizes that a shallow network with relatively small number of units works well for most of the problems. So sometimes it doesn't necessarily need to go deep and/or big.
4. **Exploration-Exploitation.** The final piece of my success is attributed to the correct implementation of an exploration strategy in the action space. I ended up applying an exponentially decayed "amplitude modulation (AM)" to the noise generation process. This AM allows the agent to explore more in the beginning of the training process and gradually converges to a reasonably well policy it has learned overtime. As a result, learning became much faster and more stable compared to the case when AM was not in place.

### Ideas for Future Work

In the future, I'm interested in exploring some other policy gradient methods such as PPO and A3C. Also, I'd like to do a hyperparameter search to find the optimal balance between performance and training time.

### Reference

- [1] T. Lillicrap, J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. *Continuous control with deep reinforcement learning*. ICLR 2016
- [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. *Human-level control through deep reinforcement learning*. Nature, 518 (7540):529–533, 2015.