

Deep Reinforcement Learning – Project 2 : Continuous Control

Melan Vijayaratnam

This document presents a technical description of the Continuous Control project in the context of the Deep Reinforcement Learning Nanodegree from Udacity.

1 Summary

For this project, we will work with the [Reacher](#) environment.

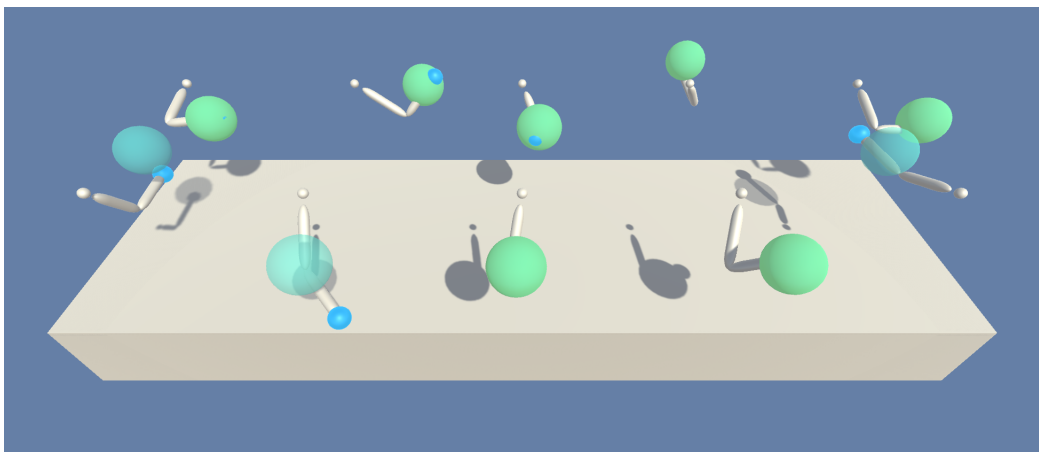


Figure 1: Unity ML-Agents Reacher environment

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 your agent is to maintain its position at the target location for as many steps as possible.

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.

This report will focus on one of the two versions of the environment, that is the one that contains a single agent. The task is **episodic** where each episode has 1000 timesteps. In order to solve the environment, the agent must get an average score of +30 over 100 consecutive episodes

2 Methods

2.1 Policy-based & value-based methods

With **value-based methods**, the agent uses its experience with the environment to maintain an estimate of the optimal action-value function. The optimal policy is then obtained from the optimal action-value function estimate:

$$\text{Interaction} \rightarrow \text{Optimal Value Function } q_* \rightarrow \text{Optimal Policy } \pi_*$$

Value-based methods

Policy-based methods directly learn the optimal policy, without having to maintain a separate value function estimate:

$$\text{Interaction} \rightarrow \text{Optimal Policy } \pi_*$$

Policy-based methods

One of the limitations of value-based methods is that they tend to a **deterministic** or **near-deterministic** policies.

On the other hand, policy-based methods can learn either **stochastic** or **deterministic** policies, so that they can be used to solve environments with either finite or continuous action spaces.

2.2 Actor-critic methods

In **actor-critic methods**, we are using value-based techniques to further reduce the variance of policy-based methods.

Basically actor-critic are a hybrid version of the policy- and value- based methods, in which the actor estimates the policy and the critic estimates the value function.

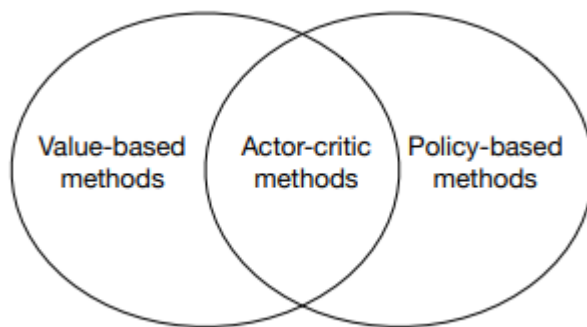


Figure 2: Actor-critic and relation to other Reinforcement Learning methods

3 Deep Deterministic Policy Gradient (DDPG)

3.1 Overview

DDPG (Lillicrap et al., 2015), short for **Deep Deterministic Policy Gradient**, is a model-free off-policy actor-critic algorithm, combining **DPG** with **DQN**. Recall that DQN (Deep Q-Network) stabilizes the learning of the Q-function by using experience replay and a frozen target network. The original DQN works in discrete space, and DDPG extends it to continuous space with the actor-critic framework while learning a deterministic policy.

In DDPG, we use 2 deep neural networks : one is the actor and the other is the critic:

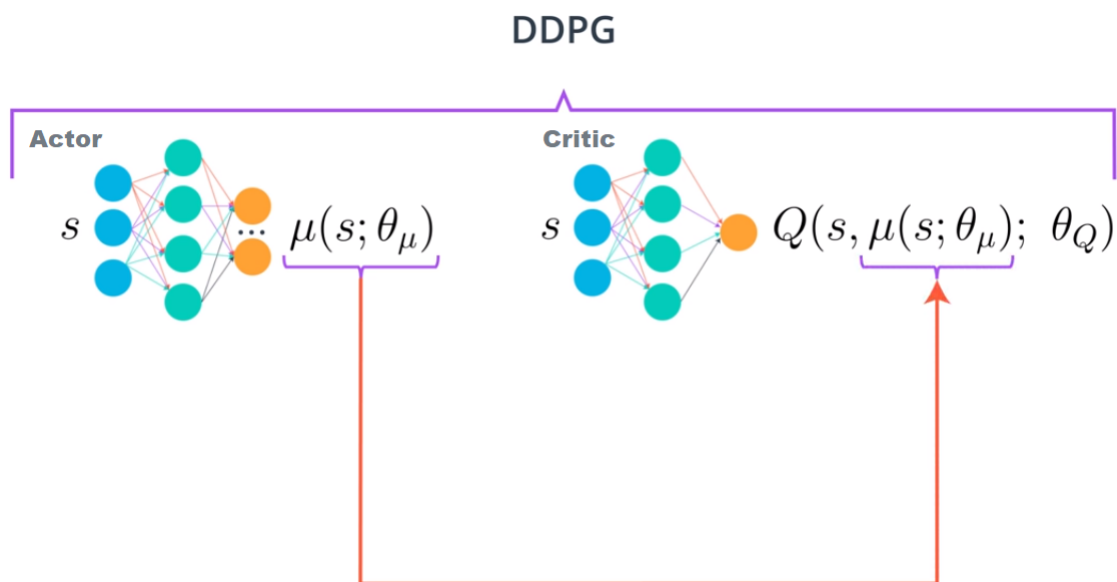


Figure 3: DDPG structure

The actor is used to approximate the optimal policy deterministically. That means we always want to output the best believed action for any given state. This is unlike a stochastic policy which learn a probability distribution over the actions.

In DDPG, we want the believed best action every single time we query the actor network, that is a deterministic policy. The actor is basically learning $\operatorname{argmax}_a Q(s, a)$ which is the best action.

The critic learns to evaluate the optimal action-value function by using the actor's best believed action.

3.2 Replay Buffer

When using neural networks, it is usually assumed that samples are independently and identically distributed (i.i.d). Alas in Reinforcement Learning, samples are generated from exploring sequentially in an environment, resulting in the previous assumption that no longer holds. Action A_t is partially responsible for the reward and state at time $(t + 1)$:

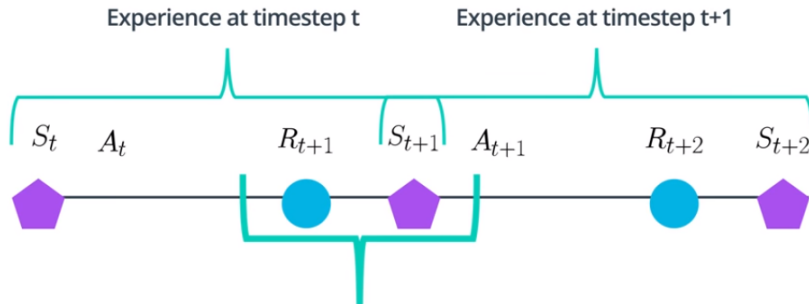
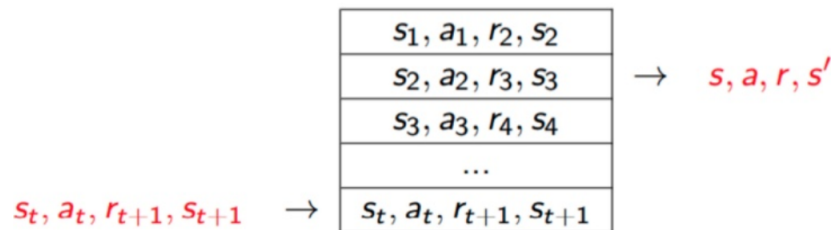


Figure 4: Highlight of correlation of sequential samples in Reinforcement Learning

As in DQN, we use a replay buffer to address this issue. Because DDPG is an off-policy algorithm, that is it employs a separate behavior policy independent of the policy being improved upon, the replay buffer can be large, allowing the algorithm to benefit from learning across a set of uncorrelated transitions.



Replay Buffer – fixed size

Figure 5: Replay buffer overview

3.3 Soft target updates

In DDPG, not only we have a regular network for the actor as well for the critic, we do also have a copy of these regular networks. We have a target actor and target critic, $Q'(s, a|\theta^{Q'})$ and $\mu'(s|\theta^{\mu'})$ respectively that are used for calculating the target values.

We update those target networks using a **soft-update strategy**. A soft-update strategy consists of slowly blending your regular network weights with your target network weights:

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

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

,with $\tau \ll 1$, usually 1% or 0.01.

This means that the target values are constrained to change slowly,
greatly improving the stability of learning.

4 Implementation details