# Udacity's Deep Reinforcement Learning Nanodegree Project 2 Report:

December 24, 2018

## 1   Introduction

## 2   Learning Algorithm

The Deep Deterministic Policy Gradient (DDPG) algorithm [1, 2] is implemented in this project:

1: Initialize replay memory $D$ with capacity $N$
2: Initialize critic network $\hat{q}$ with random weights $w^q$
3: Initialize actor network $\mu$ with weights $w^\mu$
4: Initialize target critic weights $w^{q-} \leftarrow w^q$
5: Initialize target actor weights $w^{\mu-} \leftarrow w^\mu$
6: **for** the episode $e \leftarrow 1$ to $M$ **do**
7:     Initialize a random process $\mathcal{N}$ for action exploration
8:     Receive initial input state $S$
9:     **for** time step $t \leftarrow 1$ to $T$ **do**
10:         Choose action $A = \mu(S|w^\mu) + \mathcal{N}$
11:         Take action $A$, observe reward $R$, and next input frame $S'$
12:         Store experience $(S, A, R, S')$ in replay memory $D$
13:         $S \leftarrow S'$
14:         Obtain minibatch of tuples $(s_j, a_j, r_j, s_{j+1})$ from $D$ of size $K$.
15:         Set target $y_j = r_j + \gamma \hat{q}(s_{j+1}, \mu(S', w^{\mu-}), w^{q-})$
16:         Update $w^q$: $\Delta w^q = -\alpha \frac{1}{N} \sum_j (y_j - \hat{q}(s_j, a_j, w^q)) \nabla_{w^q} \hat{q}(s_j, a_j, w^q)$
17:         Update $w^\mu$ with policy gradient:

$$\nabla_{w^\mu} J \approx \frac{1}{N} \sum_i \nabla_a \hat{q}(s_i, a\mu(s_i)) \nabla_{w^\mu} \mu(s_i|w^\mu)$$

18:         Soft update $w^{q-}$: $w^{q-} \leftarrow (1-\tau)w^{q-} + \tau w^q$
19:         Soft update $w^{\mu-}$: $w^{\mu-} \leftarrow (1-\tau)w^{\mu-} + \tau w^\mu$

# 3 Results and Plots of Rewards

# 4 Ideas for Future Work

# References

[1] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, sep 2015.

[2] Udacity. Deep Reinforcement Learning Nanodegree Course Material, 2018.