

Project 2 (Continuous control) for Udacity Deep Reinforcement Learning Nanodegree

Goal

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 time steps as possible. The following two success criteria have been considered.

- Option 1: Solve the First Version

The task is episodic, and in order to solve the environment, your 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, your agents must get an average score of +30 (over 100 consecutive episodes, and over all agents). Specifically,

After each episode, we add up the rewards that each agent received (without discounting), to get a score for each agent. This yields 20 (potentially different) scores. We then take the average of these 20 scores.

This yields an average score for each episode (where the average is over all 20 agents).

In this project, the second option is used for success criteria. The main reason is that the single agent training is slower than the 20 agents training in terms of experience collection and reduce correlations

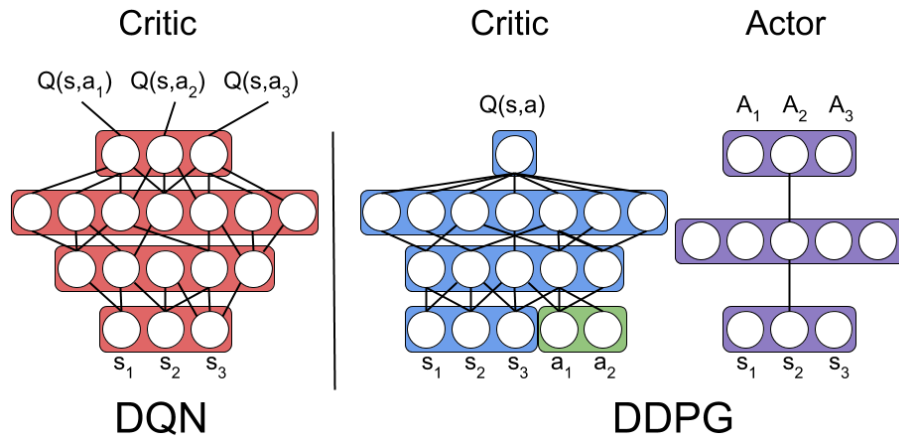
DDPG in detail

Deep Deterministic Policy Gradient (DDPG) is an algorithm which concurrently learns a Q-function and a policy. It uses off-policy data and the Bellman equation to learn the Q-function and uses the Q-function to learn the policy. A high-level DDPG structure looks the following, and you can see it has some DQN features like the replay buffer, critic network and so on. As mentioned earlier: computing the maximum over actions in the target is a challenge in continuous action spaces. DDPG deals with this by using a target policy network to compute an action which approximately maximizes $Q_{\phi_{\text{target}}}$. The target policy network is found the same way as the target Q-function: by polyak averaging the policy parameters over the course of training.

Putting it all together, Q-learning in DDPG is performed by minimizing the following MSBE loss with stochastic gradient descent:

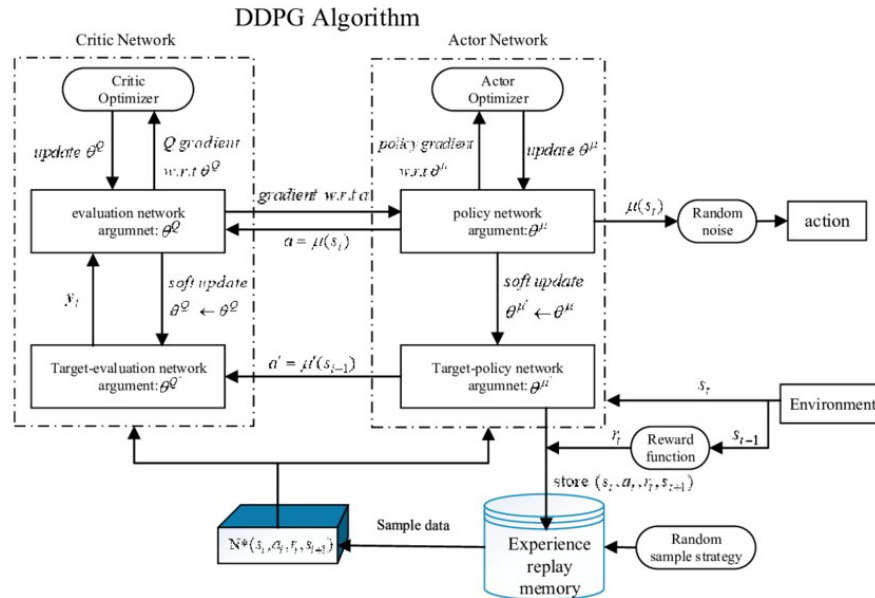
$$L(\phi, \mathcal{D}) = \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[\left(Q_{\phi}(s, a) - (r + \gamma(1-d)Q_{\phi_{\text{target}}}(s', \mu_{\theta_{\text{target}}}(s'))) \right)^2 \right],$$

The below image shows the comparison between DDPG and DQN.



Approach

The high-level structure shows as the following, and the code under ddpq_agent.py follows the diagram:



1. The state and action space of this environment

We can set up the environment (20 agents) with the following code, and the reward of +0.1 is provided for each step that the agent's hand is in the goal location. 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.

```
# reset the environment
env_info = env.reset(train_mode=True)[brain_name]

# number of agents
num_agents = len(env_info.agents)
print('Number of agents:', num_agents)
```

```

# size of each action
action_size = brain.vector_action_space_size
print('Size of each action:', action_size)

# examine the state space
states = env_info.vector_observations
state_size = states.shape[1]
print('There are {} agents. Each observes a state with length: {}'.format(states.shape[0],
state_size))
print('The state for the first agent looks like:', states[0])

```

2. Explore the environment by taking random actions

We then take some random actions based on the environment we created just now, and see how the agents perform (apparently it will be bad without learning)

```

env_info = env.reset(train_mode=False)[brain_name] # reset the environment
states = env_info.vector_observations              # get the current state (for each agent)
scores = np.zeros(num_agents)                     # initialize the score (for each agent)

while True:
    actions = np.random.randn(num_agents, action_size) # select an action (for each agent)
    actions = np.clip(actions, -1, 1)                 # all actions between -1 and 1
    env_info = env.step(actions)[brain_name]          # send all actions to the environment
    next_states = env_info.vector_observations         # get next state (for each agent)
    rewards = env_info.rewards                        # get reward (for each agent)
    dones = env_info.local_done                      # see if episode finished
    scores += env_info.rewards                        # update the score (for each agent)
    states = next_states                             # roll over states to next time step
    if np.any(dones):                                # exit loop if episode finished
        break
print('Total score (averaged over agents) this episode: {}'.format(np.mean(scores)))

```

3. Implement the DDPG algo to train the agent

The last step is to implement the DDPG algorithm to train the agents. The code can be found in `ddpg_agent.py`. However, I would like to mention several techniques to improve the speed and convergence:

- Use 20 agents to improve training speed: "Often, cooperation among multiple RL agents is much more critical: multiple agents must collaborate to complete a common goal, expedite learning, protect privacy, offer resiliency against failures and adversarial attacks, and overcome the physical limitations of a single RL agent behaving alone.", the related discussions can be found in this repo: Udacity discuss channel
- Adjust the OU noise by adding decreasing factors, and related discussions can be found in this repo: Udacity discuss channel
- Change different discount factor GAMMA to see the performance. The agent does not need to see too far to predict its next movement. So slightly reduce the GAMMA value to focus more on the current states.

Algorithm 1 DDPG algorithm

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights θ^Q and θ^μ .
Initialize target network Q' and μ' 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 $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'}))$
 Update critic by minimizing the loss: $L = \frac{1}{N} \sum_{i=1}^N (y_i - Q(s_i, a_i|\theta^Q))^2$
 Update the actor policy using the sampled policy gradient:

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

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

Results:

The average rewards along with the training process show as following:

