

Report

CRITERIA

Learning Algorithm

The report clearly describes the learning algorithm, along with the chosen hyperparameters. It also describes the model architectures for any neural networks.

1) Learning algorithm

MADDPG is used – **2 pairs of the actor and critic** (~ 1 brain, 1 external brain, **num_agents = 2**) are included in multi_agent in this project. Unlike a single-agent RL, **multiple agents**(interacting with the environment) **interact with each other**. Here, a **shared replay buffer** is implemented for the that multi agent systems(MAS). It is important to notice that many real problems should be designed by multi agents. It is known that MAS can enhance overall system performance, specifically along the dimensions of computational efficiency, reliability, extensibility, robustness, maintainability, responsiveness, flexibility, and reuse.

From now on, **most of the characteristics below** were actually also **made use of in the project 2**.

Basically, as in the project 2, here we also deal with a continuous task. This is the reason why DDPG agents are used to compose MAS. In DDPG algorithm, there still remain some important characteristics in DQN algorithm such as 1) **the critic's** implementing the network structure similar to the Q-Network in DQN and **following the same paradigm as in Q-learning** for training and 2) implementing a **replay buffer** in order to get experiences from the agent. In this project, the actor implements a current policy to try optimally mapping states to a desired action (to be used for the critic's input). With the gradients from maximizing the estimated Q-value from the critic, the actor is trained. In this way, **the critic and actor are trained interactively**.

Here, **Ornstein-Uhlenbeck noise**(to encourage exploration during training) is added to the actor's selected actions. Unlike for discrete action spaces (where exploration is done via epsilon-greedy or Boltzmann exploration, for instance), exploration is done adding noise to the action itself.

(Reference: T. Lillicrap et al., Continuous control with deep reinforcement learning (2015).)

Additionally, **Epsilon decay** is implemented to decrease an average scale of noise applied to actions.

On the actor's last output, the **Tanh activation** is used, which ensures that the range of every entry in the action vector is between -1 and 1.

For an optimizer, **Adam** is used for both actor and critic networks(actor learning rate: 1e-4, critic learning rate: 3e-4).

Using a **soft update** strategy, target networks are updated, consisting of slowly blending in your regular network weights(actually trained, most up to date) with your target network weights(used for prediction to stabilize strain).

As in the project 2, a **learning interval** is defined to speed up the learning process – performing the learning step after a specific learning interval is passed.

Lastly, I applied **dropout for critic** - `nn.Dropout(p=0.2)` – just before the last fc layer.

2) Network architectures

Vector Observation space size (per agent): 8

Number of stacked Vector Observation: 3

Vector Action space size (per agent): 2

➔ **state_size: 24 units, action_size: 2 units**

2-1) Actor Network's architecture (mapping states to actions)

Input(batch normalized, 24 units)

fc layer 1 with Batch Normalization and ReLU (input: 24 units (state_size), output: 512 units)

fc layer 2 with ReLU (input: 512 units, output 256 units)

fc layer 3 with Tanh (input: 256 units, output: 2 units (action_size))

2-2) Critic Network's architecture (mapping (state, action) pairs to Q-values)

Input(batch normalized, 24 units)

fcs layer 1 with Batch Normalization and ReLU (input: 24 units (state_size), output: 512 units)

fc layer 2 with Concatenate[(output from layer 1, action), axis=1] & ReLU

(input: 514 units(512 + action_size), output 256 units)

fc layer 3 (input: 256 units, output: 1 unit)

3) Parameters used in DDPG algorithm

Maximum episodes: 20000

Maximum steps per episode: 2000

First target average score: 6.0

4) Parameters used in DDPG Agent

```
buffer_size = 1e+6 # replay buffer size
batch_size = 256   # minibatch size
gamma = 0.99       # discount factor
tau = 1e-3         # for soft update of target parameters
lr_actor = 1e-4     # learning rate of the actor
lr_critic = 3e-4    # learning rate of the critic
ou_mu = 0.0        # Ornstein-Uhlenbeck noise parameter
ou_sigma = 0.2     # Ornstein-Uhlenbeck noise parameter
ou_theta = 0.15    # Ornstein-Uhlenbeck noise parameter
epsilon_decay = 1e-6 # Decay value for epsilon (epsilon -> epsilon - EPSILON_DECAY)
learn_every = 10    # time-step interval for learning
num_learn = 5       # number of learning with sampleing memory
```

Plot of Rewards

A plot of rewards per episode is included to illustrate that the agents get an average score of +0.5 (over 100 consecutive episodes, after taking the maximum over both agents).

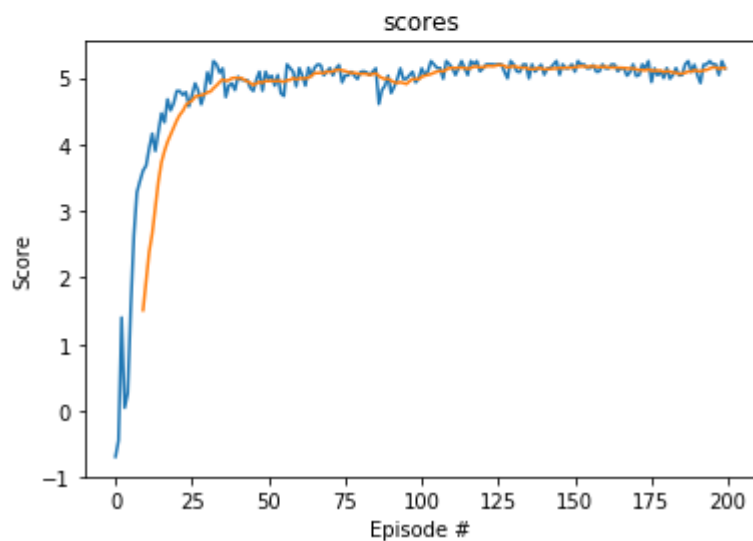
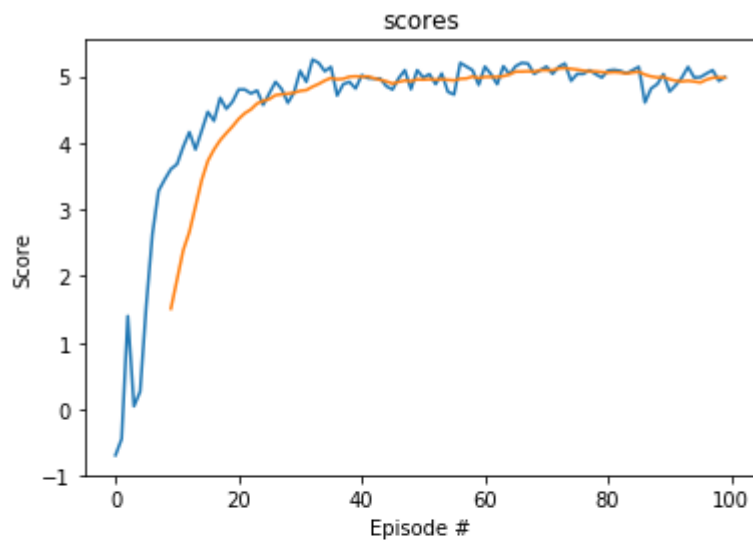
The submission reports the number of episodes needed to solve the environment.

```
scores = ddpb(target_mean_scores_deque=6.0)

# plot the scores
fig = plt.figure()
ax = fig.add_subplot(111)
plt.plot(np.arange(1, len(scores)+1), scores)
plt.ylabel('Score')
plt.xlabel('Episode #')
plt.show()

plot_scores(scores)
```

Episode 100	Average Score: 4.55
Episode 200	Average Score: 5.14
Episode 240	Average Score: 5.15



There are no problems in stopping training long after the desired score is reached. Yet, your code is built to stop training once a certain score has been reached. If I input the required score of 0.5, as I did, the training stops before one of the conditions to solve the environment is achieved, that is the average score being calculated over 100 episodes. Also, one of the requirements for the Report is to mention how many episodes it took to solve the environment. The current implementation does not give this information.

So, what to do? *Add a condition that it has been at least 100 episodes before the training is stopped.*

```

37     as train_mode:
38         for i, ddpq_agent in enumerate(multi_agent.ddpq_agents):
39             torch.save(ddpq_agent.actor_local.state_dict(), 'checkpoint_actor_' + str(i) + '.pth')
40             torch.save(ddpq_agent.critic_local.state_dict(), 'checkpoint_critic_' + str(i) + '.pth')
41
42     if np.mean(scores_deque) >= target_mean_scores_deque:
43         print('\nEnvironment solved in {} episodes! Average Score: {:.2f}'.format(i_episode, np.mean(scores_deque)))
44         if train_mode:
45             for i, ddpq_agent in enumerate(multi_agent.ddpq_agents):
46                 torch.save(ddpq_agent.actor_local.state_dict(), 'checkpoint_actor_' + str(i) + '.pth')
47                 torch.save(ddpq_agent.critic_local.state_dict(), 'checkpoint_critic_' + str(i) + '.pth')
48             break
49
50

```

Reply to the reviewer's comments:

Now the code has been modified to stop learning when `i_episode` exceeds 100 and the current average score exceeds the target score (namely 0.5).

The score above 0.5 is achieved at `i_episode=10`.

```

print('\nEpisode {} Average Score: {:.2f} (Previous Best Average Score: {:.2f})'.format(i_episode, np.mean(scores_deque), best_average_score))
if i_episode % print_every == 0:
    print('\nEpisode {} Average Score: {:.2f} (Previous Best Average Score: {:.2f})'.format(i_episode, np.mean(scores_deque), best_average_score))
    plot_scores(scores)

if np.mean(scores_deque) >= target_mean_scores_deque:
    # print only once when the environment is first resolved
    if not is_solved:
        print('\nEnvironment solved in {} episodes! Average Score: {:.2f}'.format(i_episode, np.mean(scores_deque)))
        is_solved = True
        if i_episode <= min_n_episodes:
            print('\nKeep training until i_episode is over min_n_episodes {} (current i_episode = {})'.format(min_n_episodes, i_episode))
        else:
            print('\nStop training since i_episode is over min_n_episodes {} & Average Score >= target_mean_scores_deque = {:.2f} (current i_episode = {})'.format(min_n_episodes, i_episode, np.mean(scores_deque)))
    # Save when best_average_score is renewed as well as np.mean(scores_deque) >= target_mean_scores_deque
    if train_mode and (np.mean(scores_deque) > best_average_score):
        print('\nCheckpoint files being saved...')
        for i, ddpq_agent in enumerate(multi_agent.ddpq_agents):
            print('\nSaving checkpoint_actor_' + str(i) + '.pth')
            torch.save(ddpq_agent.actor_local.state_dict(), 'checkpoint_actor_' + str(i) + '.pth')
            print('\nSaving checkpoint_critic_' + str(i) + '.pth')
            torch.save(ddpq_agent.critic_local.state_dict(), 'checkpoint_critic_' + str(i) + '.pth')

    # Keep or stop training
    if is_solved and i_episode > min_n_episodes:
        break

if np.mean(scores_deque) > best_average_score:
    print('\nBest Average Score now updated from {:.2f} to {:.2f}'.format(best_average_score, np.mean(scores_deque)))
    best_average_score = np.mean(scores_deque)

```

```

In [9]: # For reply to the reviewer's comments
# Cell for calculation to find out at which i_episode provides the first score greater than 0.5
multi_agent = MultiAgent(config, load_from_files=False)
score = ddpq(min_n_episodes=100, target_mean_scores_deque=0.5, train_mode=True)

Episode 1      Average Score: -0.69 (Previous Best Average Score: -inf)
Best Average Score now updated from -inf to -0.69
Episode 3      Average Score: -0.67 (Previous Best Average Score: -0.69)
Best Average Score now updated from -0.69 to -0.67
Episode 4      Average Score: -0.66 (Previous Best Average Score: -0.67)
Best Average Score now updated from -0.67 to -0.66
Episode 5      Average Score: -0.46 (Previous Best Average Score: -0.66)
Best Average Score now updated from -0.66 to -0.46
Episode 6      Average Score: -0.30 (Previous Best Average Score: -0.46)
Best Average Score now updated from -0.46 to -0.30
Episode 7      Average Score: -0.02 (Previous Best Average Score: -0.30)
Best Average Score now updated from -0.30 to -0.02
Episode 8      Average Score: 0.20 (Previous Best Average Score: -0.02)
Best Average Score now updated from -0.02 to 0.20
Episode 9      Average Score: 0.34 (Previous Best Average Score: 0.20)
Best Average Score now updated from 0.20 to 0.34
Episode 10     Average Score: 0.61 (Previous Best Average Score: 0.34)
Environment solved in 10 episodes!      Average Score: 0.61

Keep training until i_episode is over min_n_episodes 100 (current i_episode = 10)

Checkpoint files being saved...

Saving checkpoint_actor_0.pth

Saving checkpoint_critic_0.pth

Saving checkpoint_actor_1.pth

Saving checkpoint_critic_1.pth

Best Average Score now updated from 0.34 to 0.61
Episode 11     Average Score: 0.87 (Previous Best Average Score: 0.61)
Checkpoint files being saved...

```

Ideas for Future Work

The submission has concrete future ideas for improving the agent's performance.

- 1) Utilizing Crawler-DDPG, NAF, PPO, TRPO, TNPG, D4pG, A2C, A3C, Q-Prop, et cetera.
- 2) Implementing prioritized replay buffer.
- 3) Weighted sampling for data points with larger errors (from which the model potentially learns more)
- 4) Optimizing hyperparameters more extensively.
- 5) Utilizing Ideas in AlphaZero.
- 6) Trying other weight initialization methods.

(Choosing a good initial condition is often turned out to be very important for achieving high performance!)