

A4

Requirements

Programming language: python3

You should compare the performance of DQN and one kind of improved DQN and test them in a classical RL control environment–MountainCar. OPENAI gym provides this environment, which is implemented with python.

Introduction

In former assignments, we use state-action pair to store the values to the list. However, as the task becomes more and more complex, using deep neural network is a better choice to fit values better in high-dimensional state spaces. Such as this assignment, we need to work on a Mountain-Car problem, which has complicated states and actions.

Simply using DQN has its drawback. For example, Q-value is usually over-estimated. Therefore, we need advanced-DQN to optimize our algorithm. In this lab, I use Double DQN and Dueling DQN for comparison.

My Implementation

env

First of all, we need to introduce the environment.

```
1 import gym
2 env = gym.make('MountainCar-v0')
3 env.reset()
4 for i in range(1000):
5     env.render()
6     ...
7     env.step(env.action_space.sample()) # take a random action
8 env.close()
```

I use pytorch to implement the neural network.

DQN

For DQN, first to build the network

```
1 import torch
2 import torch.nn as nn
3
4 class Net(nn.Module):
5     def __init__(self, ):
6         super(Net, self).__init__()
7         self.fc = nn.Sequential(
```

```

8         nn.Linear(n_states, 24),
9         nn.ReLU(),
10        nn.Linear(24, 24),
11        nn.ReLU(),
12        nn.Linear(24, n_actions)
13    )
14    # optimizer and loss_func
15    self.opt = torch.optim.Adam(self.parameters(), lr=LR)
16    self.mls = nn.MSELoss()
17
18    # Given state, return action value
19    def forward(self, s):
20        return self.fc(s)

```

In DQN, first build two networks, first a eval_net, then a target_net. And then initialize a memory set to store states, actions and rewards.

I have the function to choose action by the given state and storing the parameters to the memory set and the function to learn. In the learning function, I use a delayed update to update the target network.

The code can be viewed in my project.

DDQN

As for me, DDQN and DQN is almost the same. I use the same Net class. The only difference is in the learning process. DDQN use two functions Q and Q' while updating Q

Dueling DQN

Dueling DQN changes the network structure. Now, $Q(s, a) = V(s) + A(s, a)$

```

1  class Net(nn.Module):
2      def __init__(self):
3          super(Net, self).__init__()
4          self.features = nn.Sequential(
5              nn.Linear(n_states, 16),
6              nn.ReLU(),
7              nn.Linear(16, 16),
8              nn.ReLU(),
9          )
10         self.adv = nn.Linear(16, n_actions, bias=True)
11         self.val = nn.Linear(16, 1, bias=True)
12         # optimizer and loss_func
13         self.opt = torch.optim.Adam(self.parameters(), lr=LR)
14         self.mls = nn.MSELoss()
15
16         def forward(self, x):
17             x = self.features(x)
18             adv = self.adv(x)
19             val = self.val(x)

```

```
20         x = val + adv - adv.mean()
21     return x
```

Reward Function

For a better result and a better convergence, I customize the reward function.

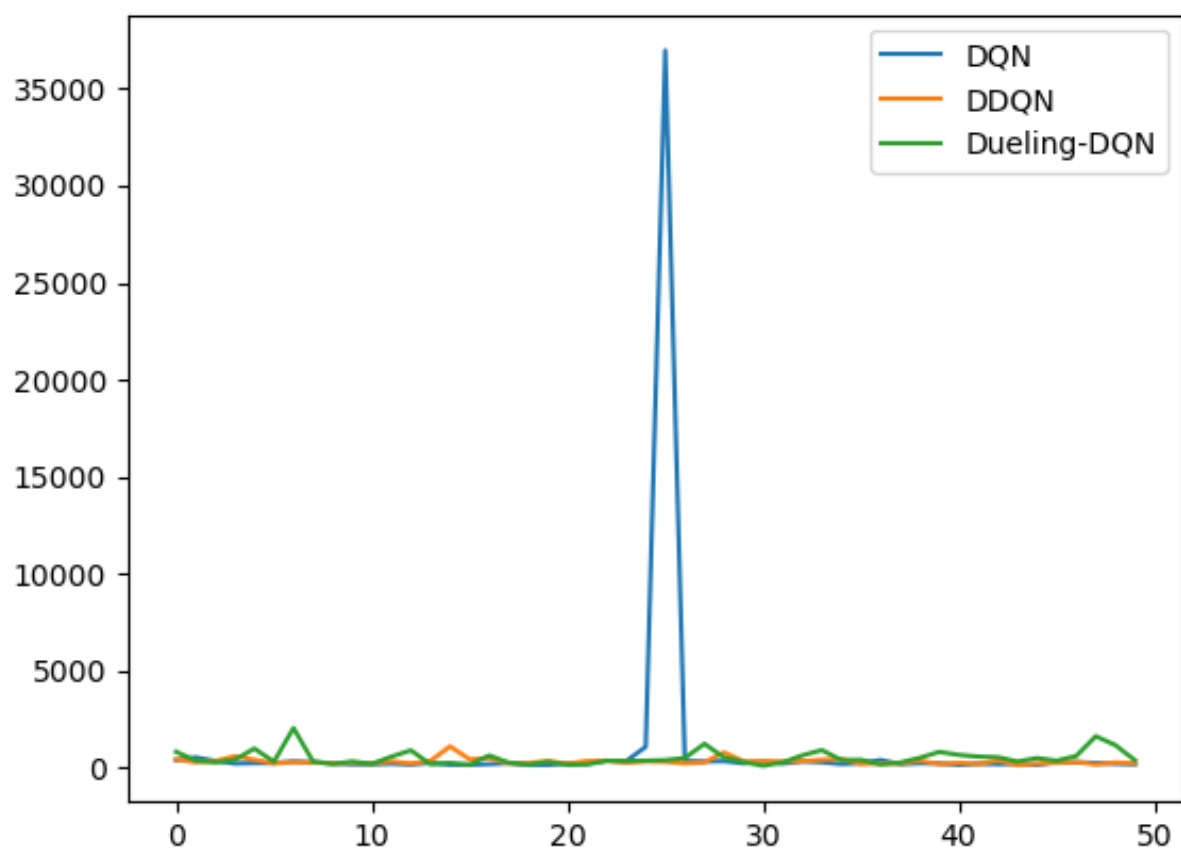
```
1  """ Rewrite the reward function, start from -0.5, reach the end
2      point at 0.5 . Because we hope the car can goes to the right side,
3      if goes lefter than start point, gives a negative reward.
4      The much closer it gets to the end point, it will receive a larger reward.
5  """
6  if -0.4 < s_[0] < 0.5:
7      r = 10 * (s_[0] + 0.4) ** 3
8  elif s_[0] >= 0.5:
9      r = 100
10 elif s_[0] <= -0.4:
11     r = -0.1
```

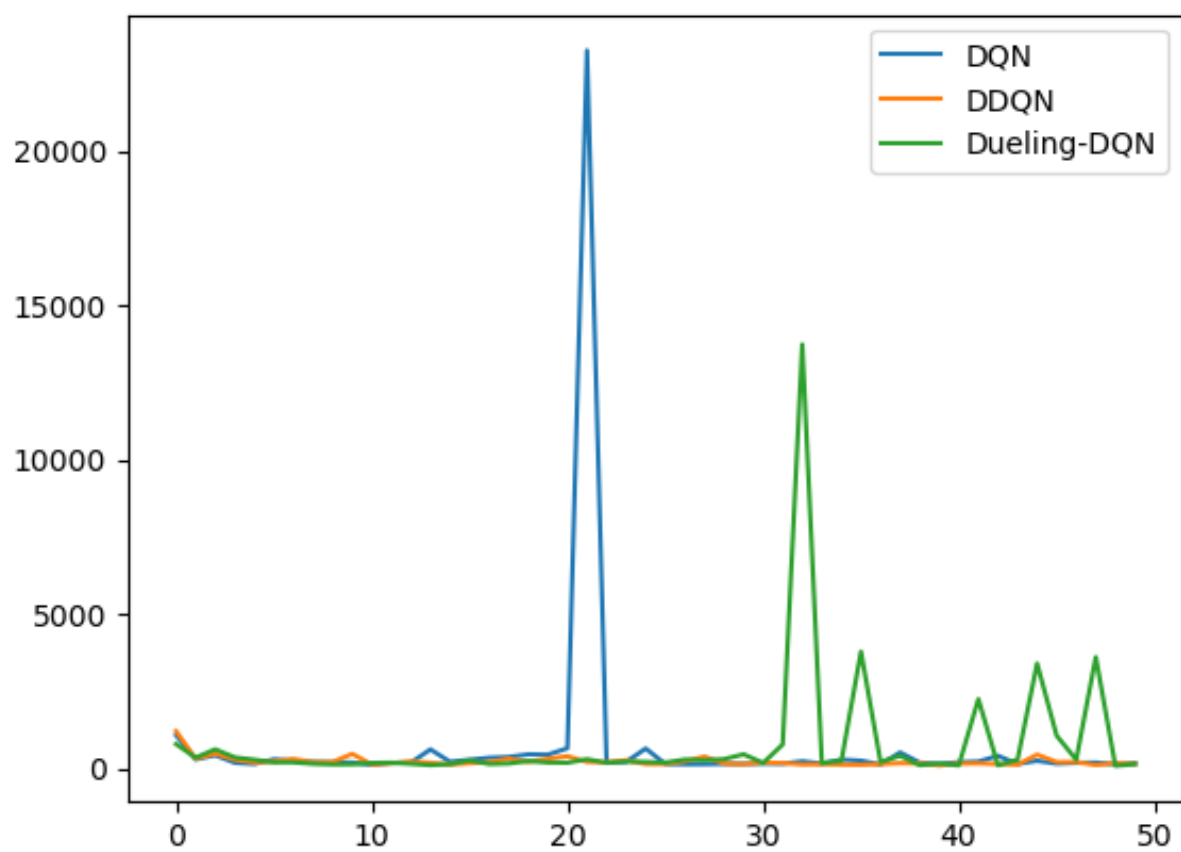
Improvements

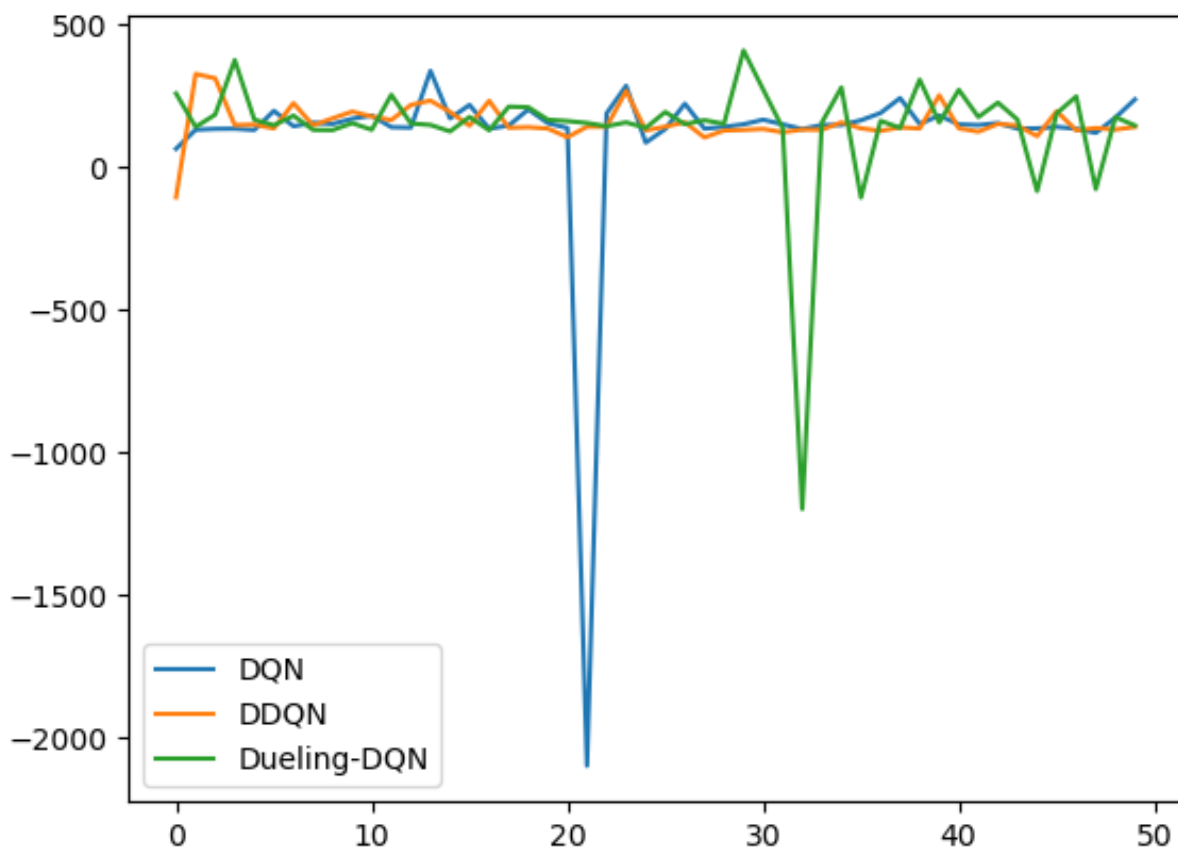
After each task, I store the current memory set to a local file. As I start the next task, the memory set will first read the information other than creating an empty set. This makes the results better, it can quickly reaches to the end point even at the first several episodes.

Result

First two pictures are **episode_num-steps** maps, and the third one is **episode_num-rewards** map.







We can see, at most times the car takes few steps(200 times for average) to reach the terminal point. The result of Double DQN is obviously more stable than DQN. Dueling-DQN also performs better than DQN.

Summary&Thinking

This lab uses pytorch and gym which I haven't use in the previous labs, which makes me a little confused when implementing it after learning the relevant knowledges. The process of learning new knowledges is also rewarding. As I do the lab, the result of DQN lets me doubt whether I have the right result. After drawing the maps, I find that I may get the right solution. Because of the time, I don't implement other advanced-DQN such as NoisyNet DQN. I will try it in the future and I wish that in the later labs I can also learn a lot.