

Assignment 4

Peng Tang 517020910038

In this assignment , we are going to do experiment with deep Q network. We are going to compare the performance of DQN, Double DQN and Dueling DQN and test them in a classical RL control environment-MountainCar.

1 Analysis

1.1 Environment Explanation

In the MountainCar game, the goal is to drive up the mountain on the right. We are going to use the following information given by the environment to train our network and reach the goal by using the trained network in the end.

- State Space: there observation is a array containing two elements. The first one is the position of the car($[-1.2, 0.6]$), and the second one is the velocity of the car($[-0.07, 0.07]$).
- Starting State: the initial position of the car is assigned a uniform random value in $[-0.6, -0.4]$. The initial velocity of the car is 0.
- Action Space: there are three actions we can take in this environment, including 0, 1, 2. 0 means accelerating to the left, 1 means not accelerating, and 2 means accelerating to the right.
- Goal: making the car's position bigger than or equal to 0.5.
- Reward: reward is 0 if the car reach the goal, -1 otherwise.
- Terminal Conditions: the car position is more than 0.5 or the episode length is greater than 200.

1.2 DQN

DQN combines the deep neural network and q-learning. We implement this algorithm on this environment as follows.

1. Initialize replay memory D to capacity N .
2. Initialize evaluation neural network Q and target neural network Q' with same weights.
3. Initialize other parameters, like episodes times M , batch size S , exploration parameter ϵ , loss function F for Q , optimizer P for Q , update interval T for updating Q' with Q every T steps, and others.
4. Initialize the environment by using the `gym` interface and get the initial observation s , then we do the following steps until this episode is over.
 - select a action a by using the ϵ -greedy algorithm
 - get the next observation s' , reward r and done d from the environment by take the action a
 - store (s, a, r, s', d) in D
 - set $s = s'$

- if D is full, sampling a batch B with size S from the D to train Q by using following steps.
 - compute evaluation value $q = Q(B(s), B(a))$
 - compute next states value $q_{next} = \max_a Q'(B(s'))$
 - compute target value $q' = B(r) + (1 - B(d)) * \lambda * q_{next}$
 - compute loss $L = F(q, q')$
 - perform a gradient descent by using L and update the weights of Q with P
- update the Q' with Q every T steps.

5. we do step 4 for M times.

By adjusting these parameters, we can get a good performance.

1.3 Double DQN

The Double DQN is used to solve the over-estimate problem of the DQN, and the difference is the way to compute the target value when training the network. Thus, we implement Double DQN as DQN and modify the step 4 in DQN as follows.

- compute evaluation value $q = Q(B(s), B(a))$
- compute the next states actions $a_{next} = \arg \max_a Q(B(s'))$
- compute next states value $q_{next} = Q'(B(s'), a_{next})$
- compute target value $q' = B(s) + (1 - B(d)) * \lambda * q_{next}$
- compute loss $L = F(q, q')$
- perform a gradient descent by using L and update the weights of Q with P

1.4 Dueling DQN

For dueling DQN, we need to change the structure of neural network. Specially, we are break the final layer in the neural network into two pieces, one is used to output $A(s, a)$ and the other is used to output $V(s)$. Then the final output of the dueling network is

$$Q(s, a) = V(s) + (A(s, a) - \frac{1}{|A|} \sum_a A(s, a))$$

2 Code

We implement the two methods in `DQN.py`.

2.1 Neural Network

```

1  class Network(nn.Module):
2      def __init__(self, input_dim, output_dim):
3          super(Network, self).__init__()
4
5          self.layer1 = nn.Sequential(
6              nn.Linear(input_dim, 32),
7              nn.ReLU(),
8              nn.Linear(32, 64),
9              nn.ReLU(),
10             nn.Linear(64, output_dim),
11
12
13         )
14

```

```

15     def forward(self, input):
16         output = self.layer1(input)
17         return output
18
19 class DuelingNet(nn.Module):
20     def __init__(self, input_dim, output_dim):
21         super(DuelingNet, self).__init__()
22
23         self.layer1 = nn.Sequential(
24             nn.Linear(input_dim, 32),
25             nn.ReLU(),
26             nn.Linear(32, 64),
27             nn.ReLU(),
28         )
29         self.layer_a = nn.Sequential(
30             nn.Linear(64, output_dim)
31         )
32         self.layer_v = nn.Sequential(
33             nn.Linear(64, 1)
34         )
35
36     def forward(self, input):
37         output1 = self.layer1(input)
38         output_a = self.layer_a(output1)
39         output_v = self.layer_v(output1)
40         output = output_v + (output_a - torch.mean(output_a))
41         return output

```

We use the `Network` class to build the Q and Q' in DQN and Double DQN and use the `DuelingNet` to build Q and Q' in Dueling DQN and Dueling&Double DQN. In the two networks, we use three linear layers and two activation function layers (`relu` function). And the `input_dim` should be 2 and the `output_dim` should be 3 in the problem.

2.2 DQN Class

This class contains two networks and one replay memory.

```

1 class DQN(object):
2     def __init__(self, input_dim, output_dim, NetworkType):
3
4         self.input_dim = input_dim
5         self.output_dim = output_dim
6         self.NetworkType = NetworkType
7
8         self.BufferSize = 10000
9         self.replayBuffer = np.zeros(self.BufferSize, dtype=object)
10        self.currentpointer = 0
11        self.batchSize = 1024
12
13        self.device = torch.device("cuda:0" if torch.cuda.is_available()
14        else "cpu")
15
16        self.evalNet = None
17        self.targetNet = None
18        self.buildNetwork()

```

```

19         self.targetNet.load_state_dict(self.evalNet.state_dict())
20         self.loss_fun = torch.nn.MSELoss().to(self.device)
21         self.optimizer = torch.optim.Adam(self.evalNet.parameters(),
lr=0.0001)
22
23         self.updateinterval = 20
24         self.epsilon = 1
25         self.gamma = 0.99
26         self.learntimes = 0
27
28     def buildNetwork(self):
29
30         #DQN or double DQN
31         if self.NetworkType == 0 or self.NetworkType == 1:
32             self.evalNet = Network(self.input_dim,
self.output_dim).to(self.device)
33             self.targetNet = Network(self.input_dim,
self.output_dim).to(self.device)
34
35         # Dueling DQN or Dueling double DQN
36         elif self.NetworkType == 2 or self.NetworkType == 3:
37             self.evalNet = DuelingNet(self.input_dim,
self.output_dim).to(self.device)
38             self.targetNet = DuelingNet(self.input_dim,
self.output_dim).to(self.device)
39
40
41     def get_action(self, obs, mode):
42
43         obs = torch.from_numpy(np.array(obs)).float().to(self.device)
44         q_val = self.targetNet(obs).detach().cpu().numpy()
45         if mode:
46             if random.random() < self.epsilon:
47                 return np.random.randint(self.output_dim)
48             else:
49                 return np.argmax(q_val)
50         else:
51             return np.argmax(q_val)
52
53
54     def push_transition(self, trans):
55         self.currentpointer = self.currentpointer % self.BufferSize
56         self.replayBuffer[self.currentpointer] = trans
57         self.currentpointer += 1
58
59     def get_traindata(self):
60         Batch = np.random.choice(self.replayBuffer, self.batchSize,
replace=False)
61         obsBatch = []
62         actionBatch = []
63         rewardBatch = []
64         nextObsBatch = []
65         doneBatch = []
66
67         for i in range(self.batchSize):
68             obsBatch.append(Batch[i][0])
69             actionBatch.append(Batch[i][1])
70             rewardBatch.append(Batch[i][2])

```

```

71         nextObsBatch.append(Batch[i][3])
72         doneBatch.append(Batch[i][4])
73
74         obsBatch = np.array(obsBatch)
75         actionBatch = np.array(actionBatch)
76         rewardBatch = np.array(rewardBatch)
77         nextObsBatch = np.array(nextObsBatch)
78         doneBatch = np.array(doneBatch)
79
80         return obsBatch, actionBatch, rewardBatch, nextObsBatch, doneBatch
81
82
83     def train(self):
84
85
86         obsBatch, actionBatch, rewardBatch, nextObsBatch, doneBatch =
self.get_traindata()
87         obsBatch = torch.from_numpy(obsBatch).float().view(-1,
2).to(self.device)
88         actionBatch = torch.LongTensor(actionBatch).view(-1,
1).to(self.device)
89         rewardBatch = torch.from_numpy(rewardBatch).view(-1,
1).to(self.device)
90         nextObsBatch = torch.from_numpy(nextObsBatch).float().view(-1,
2).to(self.device)
91         doneBatch = torch.from_numpy(doneBatch).view(-1, 1).to(self.device)
92
93         q_eval = self.evalNet(obsBatch).gather(1, actionBatch)
94         q_next = self.targetNet(nextObsBatch).detach()
95
96         # double DQN or dueling double DQN
97         if self.NetworkType == 1 or self.NetworkType == 3:
98             q_next_eval = self.evalNet(nextObsBatch).detach()
99             max_action_indx = q_next_eval.max(1)[1].view(-1,1)
100             q_target = (rewardBatch +
(torch.logical_not(doneBatch))*self.gamma*q_next.gather(1,
max_action_indx).view(-1, 1)).float()
101
102         # Dueling DQN or DQN
103         else:
104             q_target = (rewardBatch +
(torch.logical_not(doneBatch))*self.gamma*(q_next.max(1)[0]).view(-1,
1)).float()
105
106         loss = self.loss_fun(q_eval, q_target)
107         self.optimizer.zero_grad()
108         loss.backward()
109         self.optimizer.step()
110
111         if self.learntimes % self.updateinterval == 0:
112             #self.targetNet.load_state_dict(self.evalNet.state_dict())
113             for target_param, param in zip(self.targetNet.parameters(),
self.evalNet.parameters()):
114                 target_param.data.copy_(target_param.data * 0.98 +
param.data * 0.02)
115
116         self.learntimes += 1
117         self.epsilon = self.epsilon*0.99 # 下降

```

- The class implements DQN, Double DQN, Dueling DQN and Dueling&Double DQN.
- In this class, the parameters **input_dim** and **output_dim** are used to initialize the two networks; **NetworkType** is used to choose different DQNs (0 means DQN, 1 means Double DQN, 2 means Dueling DQN, 3 means Dueling&Double DQN); **replayBuffer** is used to store transitions; **BufferSize** is the size of **replayBuffer**; **currentpointer** is current position of **replayBuffer** to store transition; **batchSize** is the size of batch sampled; **device** is used to choose the gpu/cpu; **evalNet** is evaluation network and **targetNet** is target network; **loss_fun** is the loss function; **optimizer** is the optimizer to update the **evalNet**; **updateinterval** is the interval steps of updating **targetNet**; **epsilon** is the ϵ of ϵ -greedy and we will decrease it as the process; **gamma** is the λ in the equation of computing q' ; **learntimes** is memory the times of training.
- The **buildNetwork** is used to build different neural network according to the parameter **NetworkType**. If **NetworkType** is 0 or 1, we choose **Network** class; otherwise, we choose **DuelingNet** class.
- The **get_action** function is used to choose an action by giving a state. The parameter **mode** is to choose the ϵ -greedy algorithm or greedy algorithm.
- The **push_transition** function is used to store a transition in the replay memory.
- The **get_traindate** function is used to randomly sample a batch from the replay memory.
- The **train** function is used to train the evaluation network and update the target network every **updateinterval** steps. Because the only difference between the DQN and Double DQN lies in the **train** method, we can implement two different updating methods in this function to implement DQN and Double DQN. If **NetworkType** is 0 or 2, we choose DQN updating method; otherwise, we choose Double DQN updating method.

2.3 Test Function

We use this function to test the performance of our trained network by loading the parameters of our trained network. We test the score of 100 episodes to see the network's performance.

```

1  def testNetwork(parameterpath, type):
2      environment = gym.make('MountainCar-v0')
3
4      testModel = DQN(input_dim=2, output_dim=3, NetworkType=type)
5      testModel.targetNet.load_state_dict(torch.load(parameterpath))
6
7      rewardList = []
8      for i in range(100):
9          obs = environment.reset()
10         totalReward = 0
11         while True:
12             #environment.render()
13             action = testModel.get_action(obs, False)
14
15             nextobs, reward, done, info = environment.step(action)
16             totalReward += reward
17             obs = nextobs
18             if done:
19                 rewardList.append(totalReward)
20
21                 break
22         environment.close()
23

```

```
24     return rewardList
25
```

2.4 Training Function

We use this function to perform the environment and get the training data to train needed network according to given parameter.

```
1  def training(type):
2
3      environment = gym.make('MountainCar-v0')
4      maxEpisodes = 600
5      DQNModel = DQN(input_dim=2, output_dim=3, NetworkType=type)
6      step = 0
7      rewardlist = []
8      maxReward = -300
9      for episode in range(maxEpisodes):
10         observation = environment.reset()
11         totalreward = 0
12         while True:
13             step += 1
14             action = DQNModel.get_action(observation, True)
15             nextobservation, reward, done, info = environment.step(action)
16             totalreward += reward
17             reward = abs(nextobservation[0]+0.5)
18
19             transtion = (observation, action, reward, nextobservation, done)
20             DQNModel.push_transition(transtion)
21
22             observation = nextobservation
23
24             if step > DQNModel.BufferSize+1: # 存满后开始训练
25                 DQNModel.train()
26
27             if done:
28                 rewardlist.append(totalreward)
29                 print('type:', type, 'episode:', episode, 'totalreward',
totalreward)
30                 if (totalreward > maxReward and totalreward > -125) or
totalreward > -90:
31                     currentTime = time.strftime("%Y_%m_%d_%H_%M_%S",
time.localtime())
32                     savepath = 'parameters/' + str(type) + '_' + currentTime +
"_" + str(episode) + str(totalreward) + '.pkl'
33                     torch.save(DQNModel.targetNet.state_dict(), savepath)
34                     if totalreward>maxReward:
35                         maxReward = totalreward
36                     break
37
38         environment.close()
39     return rewardlist
```

In this function, we modify the reward from $r = -1$ to $r = |position + 0.5|$. In this game, the initial position is -0.5 and we need to drive up the left mountain at first to obtain a big velocity if we want to drive up the right mountain, so we give a bigger reward if we keep farer away from the initial position. In this way, we can reach our goal faster.

2.5 Main Function

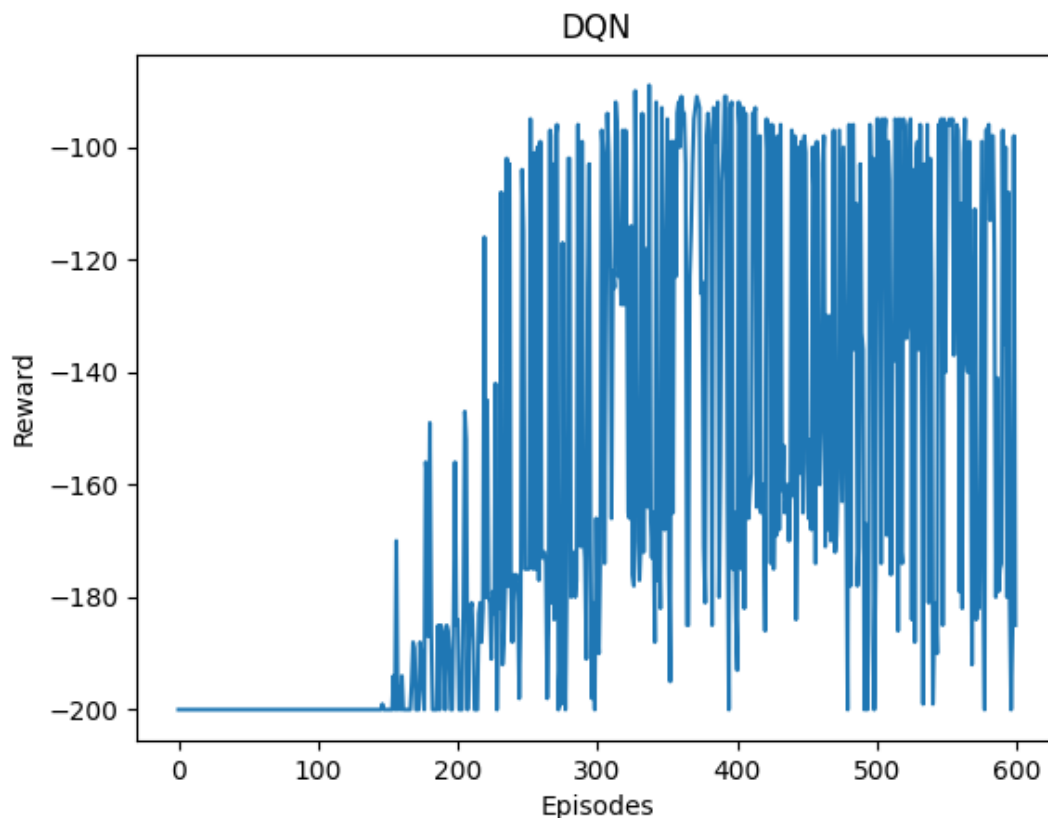
This function is used to call the **training** function by giving four different parameters and draw the final results.

```
1 def main():
2     set_seed(2)
3     DQN_Reward = training(0)
4     Double_DQN_Reward = training(1)
5     Dueling_DQN_Reward = training(2)
6     Dueling_Double_DQN_Reward = training(3)
7
8
9     drawfig(DQN_Reward, 'DQN_Reward')
10    drawfig(Double_DQN_Reward, 'Double_DQN_Reward')
11    drawfig(Dueling_DQN_Reward, 'Dueling_DQN_Reward')
12    drawfig(Dueling_Double_DQN_Reward, 'Dueling_Double_DQN_Reward')
```

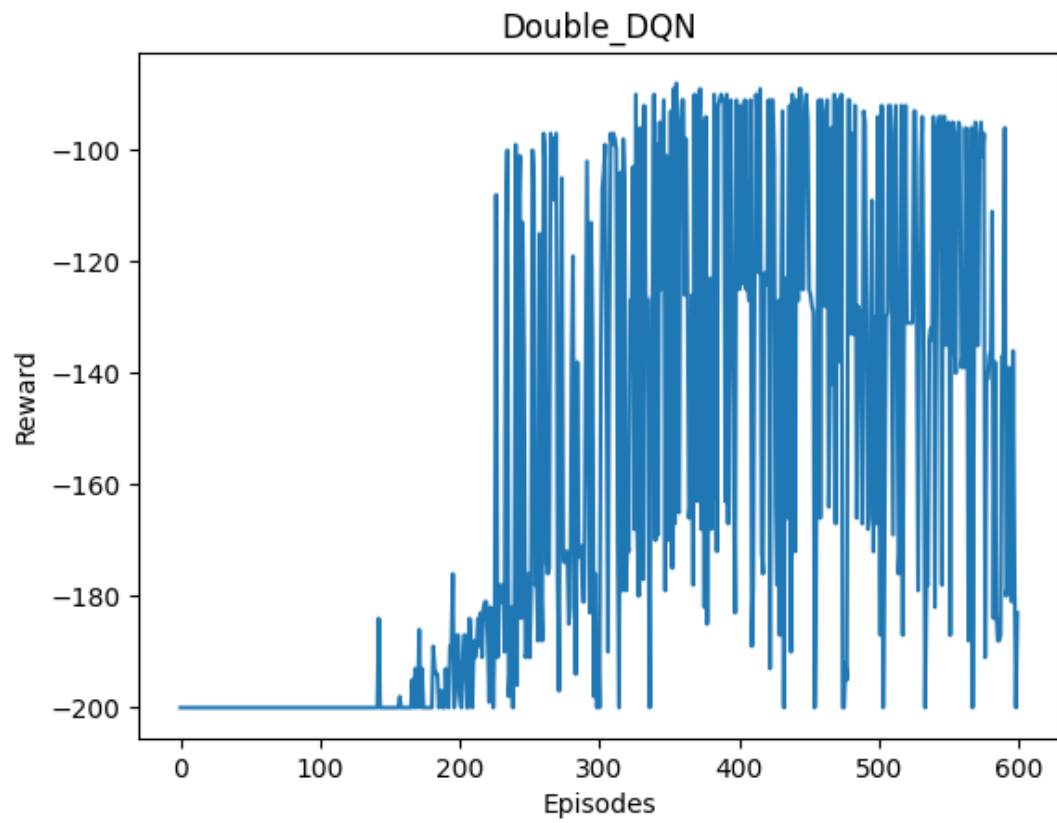
3 Experiment Results

We do experiment on the environment for four different DQNs and get the following results in training process.

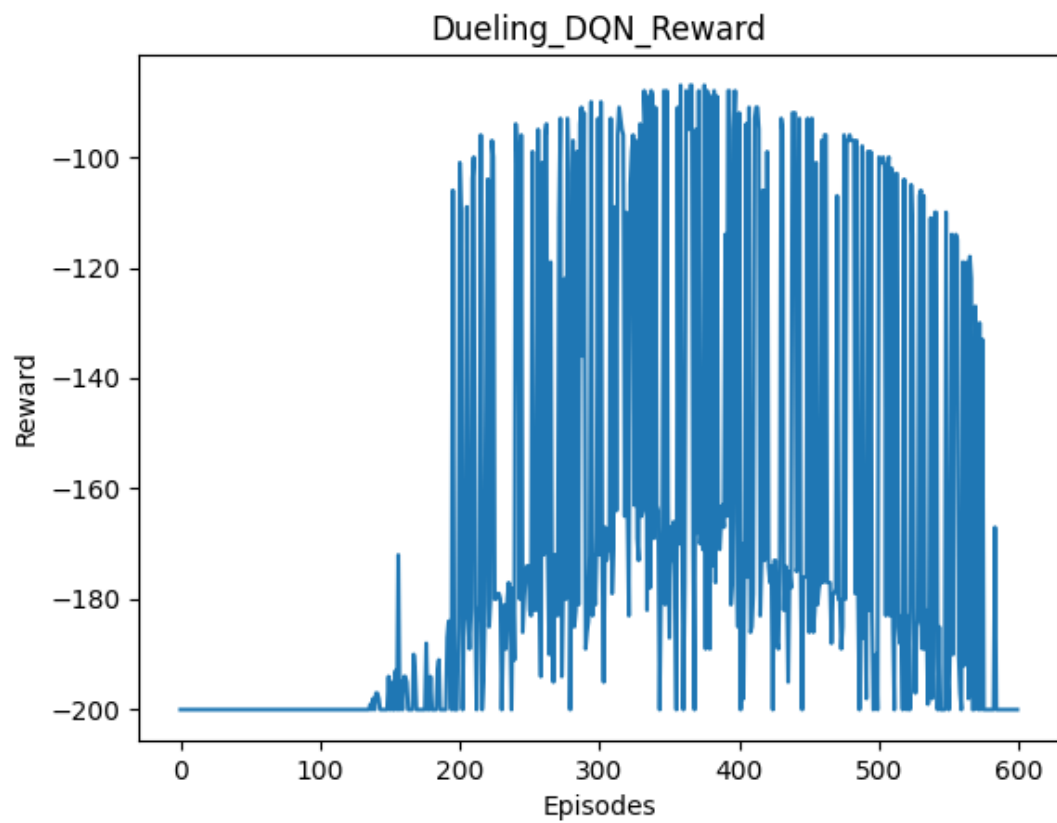
- DQN



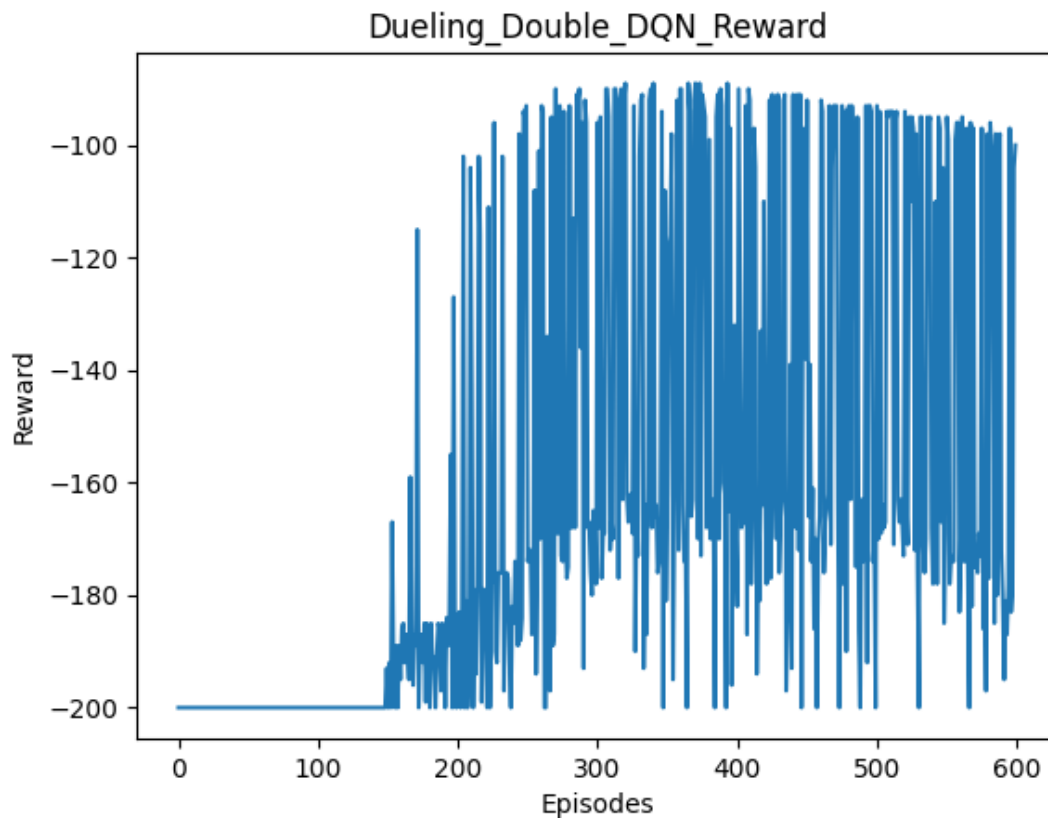
- Double DQN



- Dueling DQN

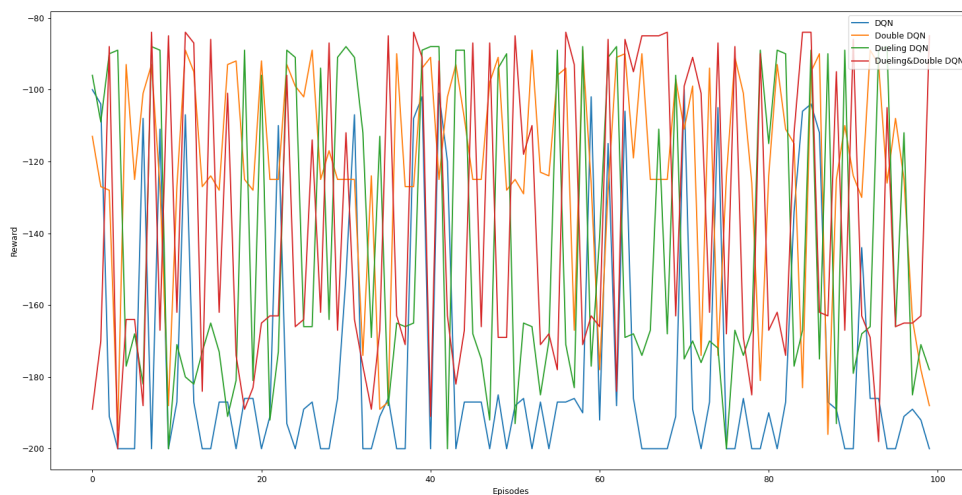


- Double&Dueling DQN



- From above figures, we can see that all methods begin to reach the goal after 150 episodes trained and all methods is fluctuation in the training process.

Then we test the trained network by performing 100 episodes on the environment and we get the following results.



```
D:\Anaconda\envs\pytorch\python.exe E:/大三下/强化学习/作业/homework4/code/DQN.py
DQN average: -175.65
Double DQN average: -121.3
Dueling DQN average: -143.11
Dueling&Double DQN average: -137.69
```

- From the two figures, we can see that Double DQN, Dueling DQN and Dueling&Double DQN have a better performance than DQN.

- Double DQN always reached the goal in the 100 episodes and Dueling DQN failed for several times.
- Dueling&Double DQN is not better than DQN in our experiment. This may be caused by hyper parameters.