

# Assignment 5

Peng Tang 517020910038

In this assignment , we are going to implement A3C and DDPG algorithm and test them in a classical RL control environment-Pendulum.

## 1 Analysis

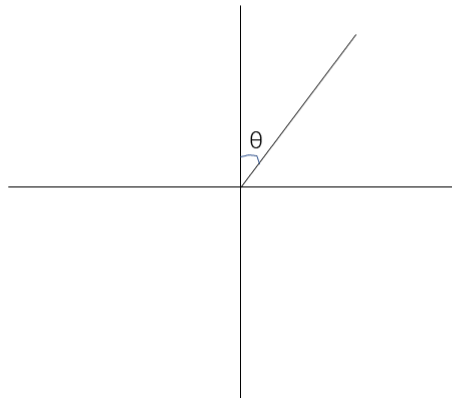
### 1.1 Environment Explanation

In the Pendulum game, the goal is to apply torque on the free end to swing the pendulum in to an upright position, with its center of gravity right above the fixed point.

- Observation Space: the observation is a array containing three elements.

Num	Observation	Min	Max
0	$x = \cos(\theta)$	-1.0	1.0
1	$y = \sin(\theta)$	-1.0	1.0
2	Angular Velocity	-8.0	8.0

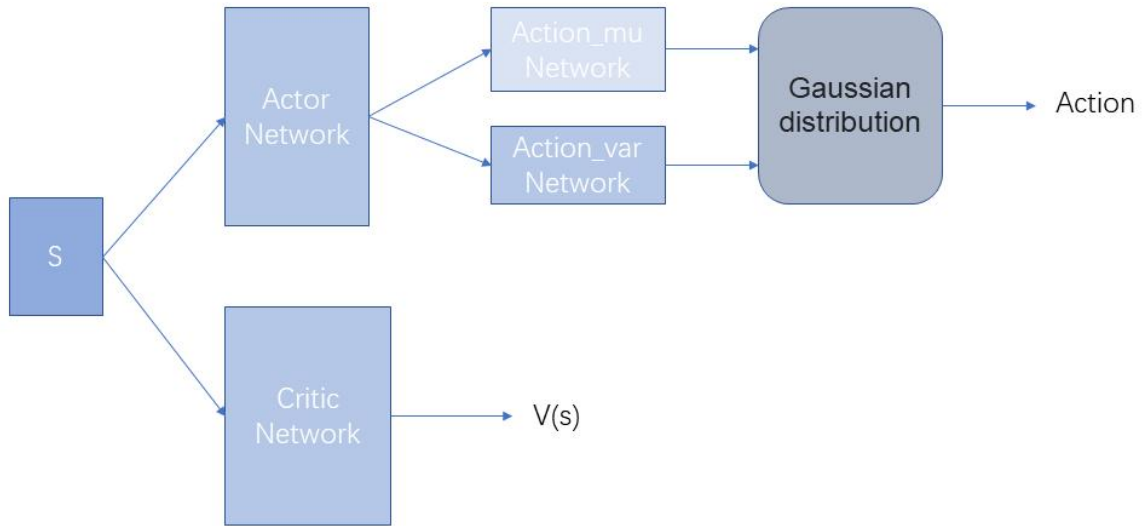
The  $\theta$  is shown as following angle.



- Starting State: The starting state is a random angle in  $[-\pi, \pi]$  and a random angular velocity in  $[-1, 1]$ .
- Action Space: The action is an array containing one element, which represents the torque applied to free end of the pendulum. It is in  $[-2, 2]$
- Reward: reward is in  $[-16.2736044, 0]$
- Terminal Conditions: The episode length is greater than 200.

### 1.2 A3C

This algorithm is a policy gradient and on-policy method. In this algorithm, we run multiple asynchronous actors-learners in parallel to speed up the training process. In specially, A3C model contains a global actor-critic network and multiple local actor-critic networks. These local actor-critic networks learn separately and share their learning results with each other through the global network. In this way, the whole model can learn more acknowledge about the environment. In the Pendulum environment, we can construct the actor-critic network as follow.



We construct Actor network to output action and Critic network to output values. Because the action space is continuous and we need to find the distribution of the action, we use the actor network to output the  $\mu, \sigma$  of action, and then choose an action from the gaussian distribution. With this network, we can implement the A3C as following steps:

- Construct a global actor-critic network  $GN$  and global counter  $T = 0$  and multiple local actor-critic network  $LN$ s with the parameters of  $GN$ .
- for each local network  $LN$ , do following steps for an episode:
  - Initialize the environment by using the `gym` interface and get the initial observation  $s$  and set a temp buffer  $N$ , the max step  $S$  and a step counter  $c = 0$ .
  - $c = c + 1$
  - select an action  $a$  by using the actor in  $LN$  with  $s$ .
  - get the next observation  $s'$ , reward  $r$ , and done  $d$  from the environment by take  $a$ .
  - store the  $(s, a, r, s', d)$  in the  $N$ .
  - if done or  $c > S$ :
    - compute the  $R_t = \sum_{i=t}^c \gamma^{i-t} r_i$  for all  $t$  by using  $r$  in  $N$  and store them in  $R$ .
    - compute the gaussian distribution  $GD$  of action and values  $V$  of states by inputing  $s$  in  $LN$ .
    - compute the actor loss  $AL = \log GD(a)(R - V)$  and critic loss  $CL = (R - V)^2$ .
    - perform a gradient descent by using  $AL$  and  $CL$  and push the gradient to  $GN$  and update the weights of  $GN$ .
    - set the weights of  $LN$  equal  $GN$ .
    - Go to the first step.

- $s = s'$ , go to the second step.

## 1.3 DDPG

DDPG is a policy gradient method and it can be considered as a continuous action version of DQN. In DQN, we use the greedy algorithm to select the action for  $s'$ . While DDPG use a network to output the action for continuous space. Especially, we use two actor-critic networks for DDPG, one is evaluation network and the other is target network. We fix the target network and train the evaluation network and update the weights of target network with the weights of evaluation for fixed steps. This is same as DQN. Because there are two networks in actor-critic network, we need to train the two networks, namely the actor and the critic. We can implement it as following steps:

- Initialize critic network  $Q$  and actor network  $A$ .
- Initialize target critic network  $Q'$  with the weights of  $Q$  and target actor network  $A'$  with the weights of  $A$ .
- Initialize the replay buffer  $N$ .
- for each episode, we do following steps:
  - Initialize the environment by using the `gym` interface and get the initial observation  $s$  and set the max step  $S$  and a step counter  $c = 0$ .
  - $c = c + 1$
  - Select an action  $a = A(s)$ .
  - Get the next observation  $s'$ , reward  $r$ , and done  $d$  from the environment by taking  $a$ .
  - store  $(s, a, r, s', d)$  in  $N$ .
  - if  $N$  is full:
    - sample a Batch  $B$  from  $N$ .
    - compute evaluation values  $q = Q(B(s), B(a))$
    - compute target values  $q' = B(r) + \gamma * Q'(B(s'), A'(B(s')))$
    - compute the loss  $CL = \frac{1}{|N|} \sum (q_i - q'_i)^2$
    - compute the actions  $a' = A(B(s))$
    - compute the values  $q'' = Q(B(s), a')$
    - compute the loss  $AL = - \sum (q''_i) / |q''|$
    - perform a gradient descent by using  $AL$  and  $CL$  and update the weights of  $A$  and  $Q$ .
  - if  $d$  is true or  $c > S$ , go the the first step
  - $s = s'$ , go to the second step.
- Until reach the max episodes.

## 2 Code

### 2.1 A3C

#### 2.1.1 A3C Network

```

1 class A3CNet(nn.Module):
2     def __init__(self, s_dim, a_dim):
3         super(A3CNet, self).__init__()
4         self.s_dim = s_dim
5         self.a_dim = a_dim

```

```

6
7     self.actorNet = nn.Sequential(
8         nn.Linear(s_dim, 256),
9         nn.ReLU(),
10        nn.Linear(256, 64),
11        nn.ReLU(),
12
13    )
14
15
16    self.actor_mu = nn.Sequential(
17        nn.Linear(64, a_dim),
18        nn.Tanh()
19    )
20    self.actor_sigma = nn.Sequential(
21        nn.Linear(64, a_dim),
22        nn.Softplus()
23    )
24
25
26    self.criticNet = nn.RNN(
27        input_size=3,
28        hidden_size=128,
29        num_layers=1,
30        batch_first=True,
31    )
32    self.out = nn.Linear(128, 1)
33
34
35
36    def forward(self, input, h_state):
37
38        a = self.actorNet(input)
39        a_mu = 2 * self.actor_mu(a)
40        a_sigma = self.actor_sigma(a) + 0.001
41
42        input1 = input[None, :, :]
43        out, h_state = self.criticNet(input1, h_state)
44        v = []
45        for i in range(out.size(1)):
46            v.append(self.out(out[:, i, :]))
47        return a_mu, a_sigma, torch.stack(v, dim=1), h_state
48
49    def choose_action(self, state, test = False):
50        h_state = None
51        action_mu, action_sigma, v, h_state =
52        self.forward(torch.FloatTensor(np.array(state)).view(-1,
53        self.s_dim), h_state)
54        h_state = h_state.data
55        if test:
56            return action_mu.data.numpy()
57        action = torch.normal(action_mu.view(1, ).data, action_sigma.view(1,
58        ).data).numpy()
59        return action
60
61    def loss_func(self, s, a, R_t):
62        h_state = None
63        self.train()

```

```

61         action_mu, action_sigma, values, h_state = self.forward(s,h_state)
62         h_state = h_state.data
63         values.view(-1,1)
64         advantages = R_t - values
65         critic_loss = advantages.pow(2)
66         distributions = Normal(action_mu, action_sigma)
67         log_prob = distributions.log_prob(a)
68         entropy = 0.5 + 0.5 * math.log(2 * math.pi) +
torch.log(distributions.scale)
69
70         actor_loss = -(log_prob * (advantages.detach())) + 0.005 * entropy)
71         return (critic_loss + actor_loss).mean()

```

- This network contains two sub networks, actor network and critic network. For actor network, we break the final layer into two sections, one for  $\mu$  and the other for  $\sigma$ . For critic network, we use RNN to implement it because RNN has a better performance on this condition.
- The **choose\_action()** is used to choose an action from the gaussian distribution that constructed by  $\mu$  and  $\sigma$ .
- The **loss\_func()** is used to compute the loss function. For convenience, we add actor loss and critic loss together.

## 2.1.2 Worker Model

```

1  class worker(mp.Process):
2      def __init__(self, name,s_dim,a_dim, gamma ,share_optimizer, globalNet,
globalmaxEpisodes, shareCurrentEpisode, shareQueue):
3          super(worker, self).__init__()
4          self.s_dim = s_dim
5          self.a_dim = a_dim
6          self.name = name
7          self.globalNet = globalNet
8          self.share_optimizer = share_optimizer
9          self.localNet = A3CNet(self.s_dim,self.a_dim)
10
11         self.globalmaxEpisodes = globalmaxEpisodes
12         self.shareCurrentEpisode = shareCurrentEpisode
13         self.shareQueue = shareQueue
14         self.gamma = gamma
15
16         self.env = gym.make('Pendulum-v1')
17
18
19     def training(self, episode_s, episode_a, episode_r, final_s, done):
20
21         h_state = None
22         R_t = []
23         s_val = 0
24         if not done:
25             action_mu, action_sigma, v, h_state
= self.localNet(torch.from_numpy(np.array(final_s)).view(-1, 3), h_state)
26             h_state = h_state.data
27
28             s_val = v.data.numpy()[0][0][0]
29
30         for r in episode_r[::-1]:

```

```

31         s_val = r + self.gamma * s_val
32         R_t.append(s_val)
33         R_t.reverse()
34
35         # 求loss
36         loss =
self.localNet.loss_func(torch.FloatTensor(np.array(episode_s)).view(-1,
self.s_dim),
37
torch.FloatTensor(np.array(episode_a)).view(-1, self.a_dim),
38
torch.FloatTensor(np.array(R_t)).view(-1, 1))
39
40         self.share_optimizer.zero_grad()
41         loss.backward()
42         for thispara1, globpara1 in zip(self.localNet.parameters(),
self.globalNet.parameters()):
43             globpara1._grad = thispara1.grad
44         self.share_optimizer.step()
45         self.localNet.load_state_dict(self.globalNet.state_dict())
46
47
48
49     def run(self):
50
51         myepisode = 0
52         while self.shareCurrentEpisode.value < self.globalmaxEpisodes:
53
54             states = []
55             actions = []
56             rewards = []
57             state = self.env.reset()
58
59             step = 0
60             totalreward = 0
61             while True:
62                 # self.env.render()
63                 step += 1
64                 action = self.localNet.choose_action(state)
65                 next_state, reward, done, _ = self.env.step(action.clip(-2,
2))
66
67                 states.append(state)
68                 actions.append(action)
69                 rewards.append((reward+8.0)/8.0)
70                 totalreward += reward
71                 if done or step >= 200:
72                     self.training(states, actions, rewards, next_state,
done)
73
74                     break
75
76                 #elif step%10 == 0:
77                 #self.training(states,actions,rewards,next_state,done)
78                 state = next_state
79
80             print(self.name, " ", myepisode, " ",
self.shareCurrentEpisode.value, " ", step, " ", totalreward)

```

```

81         if (totalreward > 0):
82             currentTime = time.strftime("%Y_%m_%d_%H_%M_%S",
time.localtime())
83             savepath1 = 'a3cparamaters/' + self.name + '_' + currentTime
+ "_" + str(self.shareCurrentEpisode.value) + "_" + str(
84                 totalreward) + '.pk1'
85             torch.save(self.localNet.state_dict(), savepath1)
86
87             myepisode += 1
88             with self.shareCurrentEpisode.get_lock():
89                 self.shareCurrentEpisode.value += 1
90                 self.shareQueue.put(totalreward)
91             self.shareQueue.put(None)
92

```

- This class is to train the network. We can define multiple objects based on this class. Then we can run them on multiple cores parallely.
- **globalNet** is the global network shared by all workers; **localNet** is the local network trained in the worker. **globalmaxEpisodes** is the number of total episodes we are going to experience.
- The **train()** is used to train the local network, update the gradient to the global network, update the weights of global network and update the local network with the global network.
- The **run()** is used to perform the environment and get the training data, then call the **train()** to train the network.

## 2.1.3 Main Function

```

1  def main():
2      env1 = gym.make('Pendulum-v1')
3      StateDim = env1.observation_space.shape[0]
4      ActionDim = env1.action_space.shape[0]
5      globalNet = A3CNet(StateDim,ActionDim)
6      globalNet.share_memory()
7      share_optimizer = SharedAdam(globalNet.parameters(),lr=0.0001)
8      MaxEpisodes = 7000
9      shareCurrentEpisode = mp.Value('i', 0)
10     shareQueue = mp.Queue()
11     gamma = 0.9
12
13     worker1 =
worker("worker1",StateDim,ActionDim,gamma,share_optimizer,globalNet,MaxEpiso
des,shareCurrentEpisode,shareQueue)
14     worker2 =
worker("worker2",StateDim,ActionDim,gamma,share_optimizer,globalNet,MaxEpiso
des,shareCurrentEpisode,shareQueue)
15     worker3 =
worker("worker3",StateDim,ActionDim,gamma,share_optimizer,globalNet,MaxEpiso
des,shareCurrentEpisode,shareQueue)
16     worker4 =
worker("worker4",StateDim,ActionDim,gamma,share_optimizer,globalNet,MaxEpiso
des,shareCurrentEpisode,shareQueue)
17
18     worker4.start()
19     worker3.start()

```

```

20     worker2.start()
21     worker1.start()
22
23     rewardlist = []
24     while True:
25         r = shareQueue.get()
26         if r is not None:
27             rewardlist.append(r)
28         else:
29             break
30     worker4.join()
31     worker3.join()
32     worker2.join()
33     worker1.join()
34
35     re = 0
36     rewardlist1 = []
37     for item in rewardlist:
38         if re == 0:
39             re = item
40         else:
41             re = re*0.95+item*0.05
42         rewardlist1.append(re)
43
44     drawfig(rewardlist1)

```

The function is used to run four processes to train the global network.

## 2.2 DDPG

### 2.2.1 Actor-Critic Network

```

1  class ActorNet(nn.Module):
2      def __init__(self, s_dim, a_dim):
3          super(ActorNet, self).__init__()
4          self.actor = nn.Sequential(
5              nn.Linear(s_dim, 128),
6              nn.ReLU(),
7              nn.Linear(128, 64),
8              nn.ReLU(),
9              nn.Linear(64, a_dim)
10         )
11         self.actor_mu = nn.Tanh()
12
13
14     def forward(self, input):
15         a = self.actor(input)
16         a_mu = 2 * self.actor_mu(a)
17         return a_mu
18
19 class CriticNet(nn.Module):
20     def __init__(self, s_dim):
21         super(CriticNet, self).__init__()
22         self.critic = nn.Sequential(
23             nn.Linear(s_dim+1, 128),
24             nn.ReLU(),

```



```

25         nn.Linear(128, 64),
26         nn.ReLU(),
27         nn.Linear(64, 1)
28     )
29     def forward(self, input):
30         return self.critic(input)

```

- The **ActorNet** is used to build  $A$ , which maps the state to the action.
- The **CriticNet** is used to build  $Q$ , which compute the q-value by inputing the  $s, a$ .
- In this game, the output layer is one dim for both networks.

## 2.2.2 DDPG Network

```

1  class DDPGNet(object):
2      def __init__(self, s_dim, a_dim):
3          self.a_dim = a_dim
4          self.s_dim = s_dim
5          self.gamma = 0.9
6
7          self.critic = CriticNet(s_dim)
8          self.critic_target = CriticNet(s_dim)
9
10         self.actor = ActorNet(s_dim,a_dim)
11         self.actor_target = ActorNet(s_dim,a_dim)
12
13         self.actor_optimizer = torch.optim.Adam(self.actor.parameters(),
14lr=0.001)
15         self.critic_optimizer = torch.optim.Adam(self.critic.parameters(),
16lr=0.001)
17
18         self.critic_target.load_state_dict(self.critic.state_dict())
19         self.actor_target.load_state_dict(self.actor.state_dict())
20         self.action_var = 0.5
21
22         self.bathsize = 256
23         self.BufferSize = 5000
24         self.currentpointer = 0
25         self.replayBuffer = np.zeros(self.BufferSize,dtype=object)
26
27         self.loss_fun = torch.nn.MSELoss()
28
29     def get_action(self, state, noise = False):
30
31         action_mu =
32         self.actor(torch.FloatTensor(np.array(state)).view(-1,self.s_dim))
33         action = torch.normal(action_mu.view(1,).data,
34         torch.from_numpy(np.array([self.action_var]))).numpy()
35
36         if noise:
37             return action_mu.view(1,).data.numpy()
38         else:
39             self.action_var= self.action_var*0.99
40             return action
41
42     def push_transition(self,trans):

```

```

41         self.currentpointer = self.currentpointer % self.BufferSize
42         self.replayBuffer[self.currentpointer] = trans
43         self.currentpointer += 1
44
45     def get_traindata(self):
46         Batch =
np.random.choice(self.replayBuffer, self.bathsize, replace=False)
47         obsBatch = []
48         actionBatch = []
49         rewardBatch = []
50         nextObsBatch = []
51         doneBatch = []
52
53         for i in range(self.bathsize):
54             obsBatch.append(Batch[i][0])
55             actionBatch.append(Batch[i][1])
56             rewardBatch.append(Batch[i][2])
57             nextObsBatch.append(Batch[i][3])
58             doneBatch.append(Batch[i][4])
59
60         obsBatch = np.array(obsBatch)
61         actionBatch = np.array(actionBatch)
62         rewardBatch = np.array(rewardBatch)
63         nextObsBatch = np.array(nextObsBatch)
64         doneBatch = np.array(doneBatch)
65
66         return obsBatch, actionBatch, rewardBatch, nextObsBatch, doneBatch
67
68
69     def train(self):
70         obsBatch, actionBatch, rewardBatch, nextObsBatch, doneBatch =
self.get_traindata()
71         obsBatch = torch.FloatTensor(obsBatch).view(-1, self.s_dim)
72         actionBatch = torch.FloatTensor(actionBatch).view(-1, 1)
73         rewardBatch = torch.FloatTensor(rewardBatch).view(-1, 1)
74         nextObsBatch = torch.FloatTensor(nextObsBatch).view(-1, self.s_dim)
75         doneBatch = torch.from_numpy(doneBatch).view(-1, 1)
76
77
78         q_eval = self.critic(torch.cat((obsBatch, actionBatch), 1))
79         nextActionBatch = self.actor_target(nextObsBatch)
80         q_next = self.critic_target(torch.cat((nextObsBatch,
nextActionBatch), 1)).detach()
81         q_target = rewardBatch +
(torch.logical_not(doneBatch)) * self.gamma * q_next
82         loss = self.loss_fun(q_eval, q_target)
83         self.critic_optimizer.zero_grad()
84         loss.backward()
85         self.critic_optimizer.step()
86
87         actionpre = self.actor(obsBatch)
88
89         q_pre = self.critic(torch.cat((obsBatch, actionpre), 1))
90         loss1 = -q_pre.mean()
91         self.actor_optimizer.zero_grad()
92         loss1.backward()
93         self.actor_optimizer.step()
94

```

```

95         for target_critic_param, critic_param in
zip(self.critic_target.parameters(), self.critic.parameters()):
96             target_critic_param.data.copy_(target_critic_param.data*0.9 +
critic_param.data*0.1)
97
98         for target_actor_param, actor_param in
zip(self.actor_target.parameters(), self.actor.parameters()):
99             target_actor_param.data.copy_(target_actor_param.data*0.9+actor_param.data
*0.1)
100

```

- In this class, we build two critic networks(**critic** and **critic\_target**) and two actor networks(**actor** and **actor\_target**). And we use the **replayBuffer** to store the transitions. The **actor\_optimizer** is for **actor** and the **critic\_optimizer** is for **critic**. The **loss\_fun** is for **critic**.
- The **get\_action()** is used to choose a action. In order to implement exploration, we use the gaussian distribution to choose a action. And as processing, we will decrease the var of the distribution, then we can use the output of  $A$  more.
- The **push\_transition()** is used to store a transition in the replay memory.
- The **get\_traindate()** is used to randomly sample a batch from the replay memory.
- The **train()** is used to train the **critic** and **actor** according to the analysis.

## 2.2.3 Running function

```

1  def running():
2      env = gym.make('Pendulum-v1')
3      StateDim = env.observation_space.shape[0]
4      ActionDim = env.action_space.shape[0]
5      maxEpisodes = 2000
6      maxstep = 400
7      DDPG = DDPGNet(StateDim, ActionDim)
8      totalstep = 0
9      rewardlist = []
10
11     for episode in range(maxEpisodes):
12         obs = env.reset()
13         totalreward = 0
14         step = 0
15         while True:
16             #env.render()
17             totalstep +=1
18             step += 1
19             action = DDPG.get_action(obs)
20             nextobs, reward, done, _ = env.step(action)
21             totalreward += reward
22             transtion = (obs,action,reward,nextobs,done)
23             DDPG.push_transition(transtion)
24
25             if totalstep>DDPG.BufferSize+1:
26                 DDPG.train()
27
28             if done or step>maxstep:
29
30                 print(episode, " ", totalreward)
31

```

```

32         if(totalreward>-125):
33             currentTime = time.strftime("%Y_%m_%d_%H_%M_%S",
time.localtime())
34             savepath1 = 'ddpgparamaters/' + "actor" + '_' +
currentTime + "_" + str(episode) + "_" + str(
35                 totalreward) + '.pkl'
36             savepath2 = 'ddpgparamaters/' + "critic" + '_' +
currentTime + "_" + str(episode) + "_" + str(
37                 totalreward) + '.pkl'
38             torch.save(DDPG.actor.state_dict(), savepath1)
39             torch.save(DDPG.critic.state_dict(),savepath2)
40
41
42             break
43             obs = nextobs
44             rewardlist.append(totalreward)
45
46         drawfig(rewardlist)

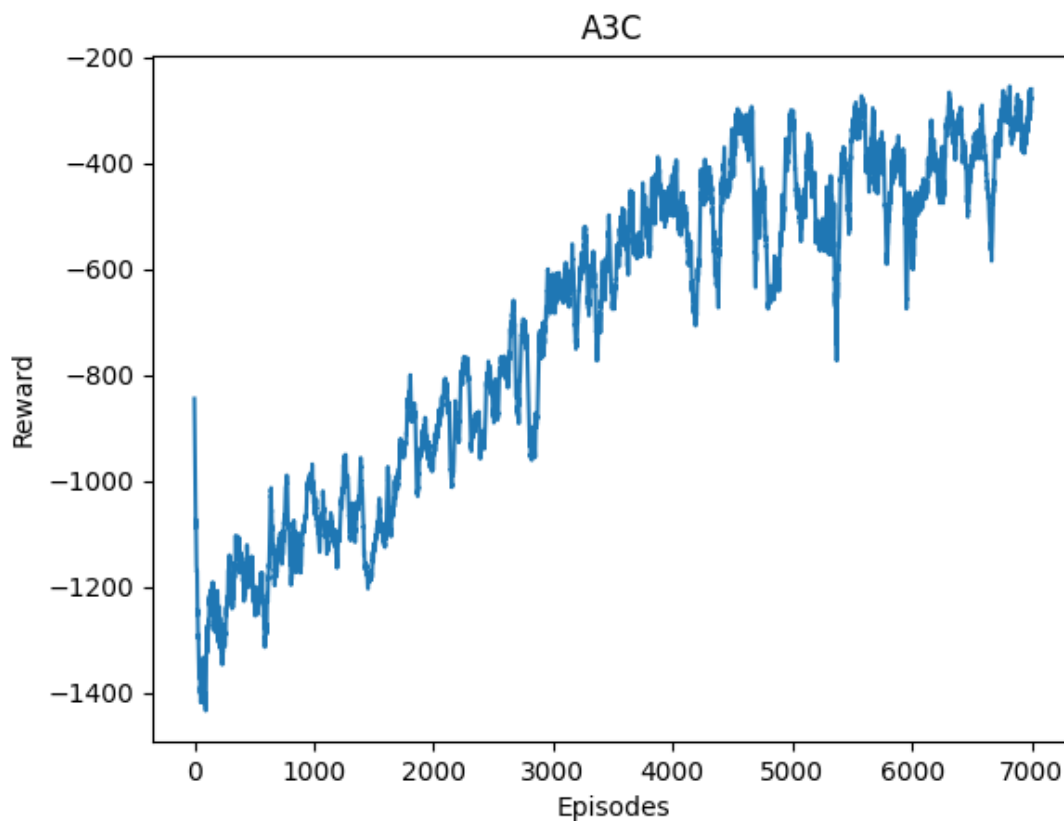
```

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

## 3 Experiment Results

### 3.1 A3C

By running the *A3C.py*, we got the following result.



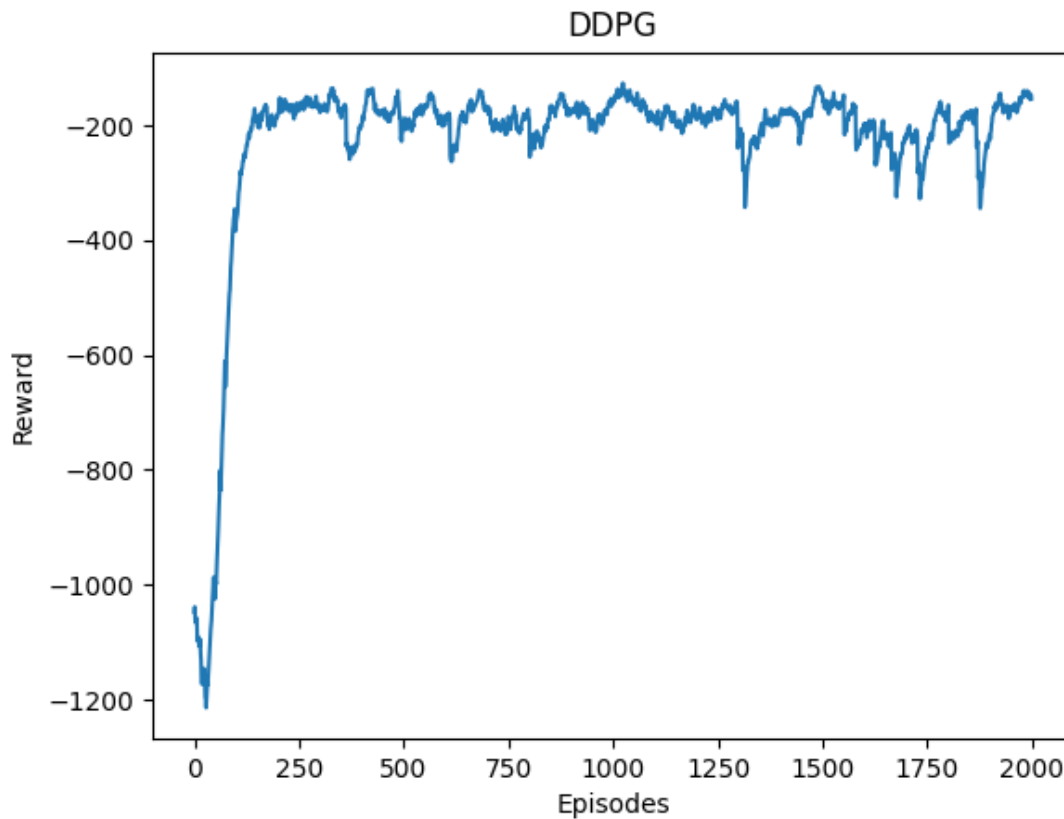
From the figure, we can see that we need 5000 episodes to reach a relatively good performance.

We tested the trained network for 100 episodes and got the following result.

```
D:\Anaconda\envs\pytorch\python.exe E:/大三下/强化学习/作业/homework5/517020910038-唐鹏/A3C.py
the average reward of 100 episodes: [-324.05457]
```

## 3.2 DDPG

By running the *DDPG.py*, we got the following result.



From the figure, we can see that, after about 250 episodes, we can get a good performance. Thus, We don't need to train it too much times in the DDPG.

We tested the trained network for 100 episodes and got the following result.

```
D:\Anaconda\envs\pytorch\python.exe E:/大三下/强化学习/作业/homework5/DDPG.py
the average reward of 100 episodes: -160.19685255953718
```

## 3.3 Compare

- From above results, we can see that the DDPG has a better performance than A3C in this environment. This is because DDPG is more suitable for continuous space problem, while A3C is more suitable for discrete space problem.
- What's more, A3C needs more episodes to train the network. But A3C runs faster than DDPG for the same number of episodes.

