

H1 CS489 Assignment 5 Report

517030910214 Hongzhou Liu

H2 0. Introduction

In this assignment, I implemented A3C in Pendulum environment. In the Pendulum environment, our aim is to keep a frictionless pendulum standing up. The only way to control the pendulum is to apply force on the joint. Thus, the action space consists of only joint effort which is in the range of $[-2.0, 2.0]$. The observation consists of the sine and cosine value of current angle θ and angular velocity $\dot{\theta}$. The reward function is

$$-(\theta^2 + 0.1 \times \dot{\theta}^2 + 0.001 \times action^2) \quad (1)$$

It reflects our goal, which is to keep the pendulum standing up ($\theta = 0$) with least angular velocity and joint effort. In this environment, the termination is not specified and I set the maximum number of steps as 200.

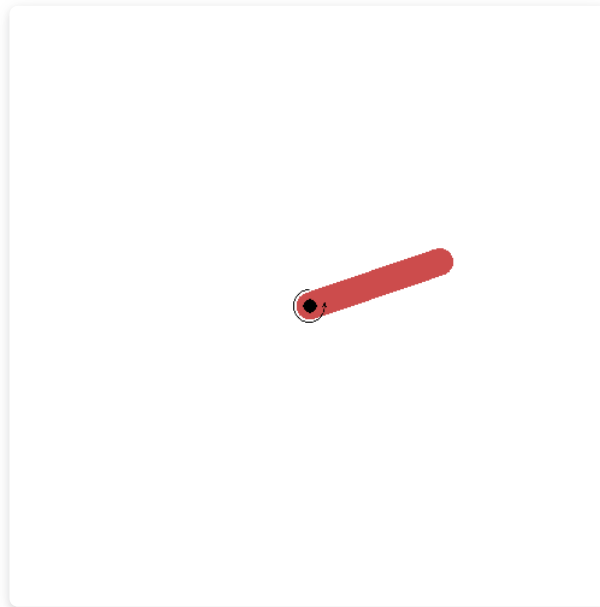


Fig.1 Pendulum

Environment:

- Ubuntu 18.04 LTS
- Python 3.7.7

H2 1. A3C

H3 1.1 Algorithm

Asynchronous Advantage Actor-Critic (A3C) is a improved algorithm of original Actor-Critic algorithm which can be used to train an agent with continuous action space. A3C consists of 3 basic ideas:

- Actor-Critic

AC algorithms take the advantages of both Q-learning and Policy Gradient. The network estimates both value function $V(s)$ and policy $\pi(s)$. The value function is seen as critic and the policy is the actor.

- Advantage

In Policy Gradient algorithm, the update rule used the discounted returns from an episode in order to tell the agent which of its actions are good and which are bad. The network will be updated in order to encourage and discourage actions appropriately.

- Asynchronous

There are multiple agents called workers who interact with environment. Each worker has its own network and parameters. There is a global network and each work will update the global network asynchronously. Such method works better than having a single agent (A2C) due to the diversity of experiences.

The algorithm will first initialize the global network and start all workers. A worker will reset gradients $d\theta, d\theta_v = 0$ at the beginning of each episode and then interacts with environment and record states, actions and rewards. After a certain steps, it will calculate the return values of each states then calculate value and policy losses and optimize the global network using the accumulated gradient of both losses by

$$d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i | s_i; \theta')(R - V(s_i; \theta'_v)) \quad (2)$$

$$d\theta_v \leftarrow d\theta_v + \partial(R - V(s_i; \theta'_v))^2 / \partial \theta'_v \quad (3)$$

where θ' and θ'_v are local parameters. Then, the algorithm will perform asynchronous update of global parameters θ and θ_v using $d\theta$ and $d\theta_v$ respectively. The algorithm will terminate after a number of steps or episodes.

H3 1.2 Implementation

H4 1.2.1 Network

We need to separate the output of our network into 2 parts, one for actor, the other for critic. As we can see in the following figure, the two outputs may share some parameters.

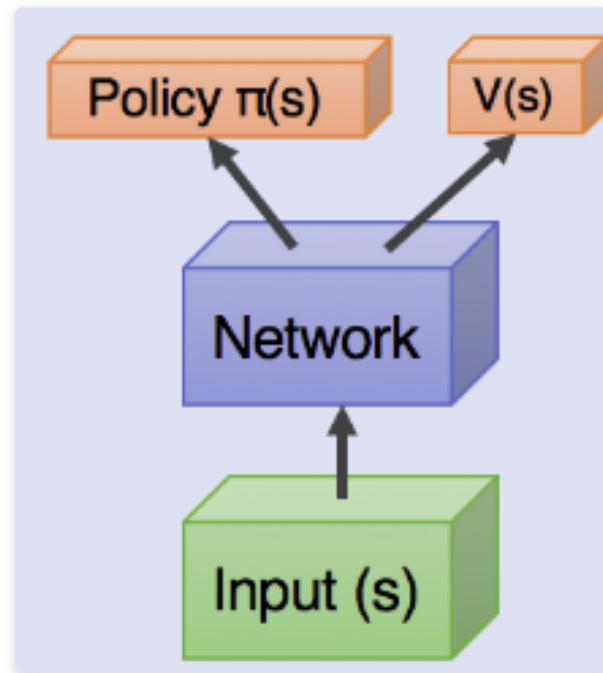


Fig.2 Network

In my network, they are totally divided except for the input layer, which mean there are an actor network and a critic network.

```
1 class ACNet(nn.Module):
```

```

2  def __init__(self, inputSize, hiddenSize, outputSize):
3      super(ACNet, self).__init__()
4      self.actor1 = nn.Linear(inputSize, hiddenSize)
5      self.mu = nn.Linear(hiddenSize, outputSize)
6      self.sigma = nn.Linear(hiddenSize, outputSize)
7
8      self.critic1 = nn.Linear(inputSize, hiddenSize)
9      self.value = nn.Linear(hiddenSize, 1)
10
11     for l in [self.actor1, self.mu, self.sigma, self.critic1, self.value]:
12         self._initLayer(l)
13
14     self.distribution = torch.distributions.Normal
15
16     def _initLayer(self, layer):
17         nn.init.normal_(layer.weight, mean=0.0, std=0.1)
18         nn.init.constant_(layer.bias, 0.0)
19
20     def forward(self, x):
21         actor1 = F.relu6(self.actor1(x))
22         mu = 2 * torch.tanh(self.mu(actor1))
23         sigma = F.softplus(self.sigma(actor1)) + 0.001
24         critic1 = F.relu6(self.critic1(x))
25         value = self.value(critic1)
26         return mu, sigma, value

```

The policy $\pi(s)$ is actually a probabilistic distribution. Here, we assume it's a Gaussian distribution and the network learns it's mean and variance. As for the critic $V(s)$, we use a network with single hidden layer and outputs it's value directly. The choice of activation functions refers to [1](#). I also defined loss function like

```

1  def loss(self, state, action, R):
2      self.train()
3      mu, sigma, value = self.forward(state)
4      error = R - value
5      critic_loss = error.pow(2)
6      dist = self.distribution(mu, sigma)
7      log_prob = dist.log_prob(action)
8      entropy = 0.5 + 0.5 * math.log(2 * math.pi) + torch.log(dist.scale)
9      actor_loss = -(log_prob * error.detach() + 0.005 * entropy)
10     return (critic_loss + actor_loss).mean()

```

It's almost the same as the original loss function in the paper. However, I added an entropy term in actor loss to enable a few explorations.

We also need to select actions during training. It's quite difficult from what we do in DQN. Here, the policy is probabilistic instead of a deterministic one. Thus, we need to sample from the distribution generated from the outputs of the neural network.:

```

1  def selectAction(self, state):
2      self.training = False
3      with torch.no_grad():
4          mu, sigma, _ = self.forward(state)
5          dist = self.distribution(mu.detach(), sigma.detach())
6          return dist.sample().numpy()

```

To train the global network and local networks, we also need to modify the optimization method.

```

1  class SharedAdam(optim.Adam):
2      def __init__(self, params, lr=1e-3, betas=(0.9, 0.99), eps=1e-8,
3          weight_decay=0):
4          super(SharedAdam, self).__init__(params, lr=lr, betas=betas, eps=eps,
5              weight_decay=weight_decay)
6          for group in self.param_groups:
7              for p in group['params']:
8                  state = self.state[p]
9                  state['step'] = 0
10                 state['exp_avg'] = torch.zeros_like(p.data)
11                 state['exp_avg_sq'] = torch.zeros_like(p.data)
12
13                 state['exp_avg'].share_memory_()
14                 state['exp_avg_sq'].share_memory_()

```

The `SharedAdam` make some of its parameters shared thus can be used in A3C.

H4 1.2.2 Parallel Learning

We need to utilize `multiprocessing` provided by `torch` to make a bunch of workers learn while interacting with the environment. A single process is implemented as follows.

```

1  class Worker(mp.Process):
2      def __init__(self, rank, globalNet, localNet, optimizer, totalEpisode,
3          globalReturn, Q, params):
4          super(Worker, self).__init__()
5          self.rank = rank
6          self.env = gym.make('Pendulum-v0').unwrapped
7          self.GNet = globalNet
8          self.LNet = localNet
9          self.opt = optimizer
10         self.totEps = totalEpisode
11         self.totR = globalReturn
12         self.Q = Q
13         self.params = params

```

Here, `LNet` is the local network of this work. `totEps` and `totR` are variables resides in the shared memory space. We can use them to record some information among workers. `Q` is a queue which can be seen as another segment of memory that some workers can put data into it and other workers can get those data from it. `params` are some hyper parameters.

The training process is implemented in `run` function that will be executed after the parent process called `workerProcess.start()`. And our A3C algorithm is implemented in this part.

```
1  def run(self):
2      steps = 1
3      while self.totEps.value < self.params['MAX_EPISODE']:
4          stateBuf, actionBuf, rewardBuf = [], [], []
5          state = self.env.reset()
6          ret = 0.0
7          rewardDecay = 1.0
8
9          for t in range(self.params['MAX_STEP']):
10             if self.rank == 0:
11                 self.env.render()
12                 action = self.LNet.selectAction(torch.from_numpy(state.reshape(1,
-1).astype(np.float32)).to(device))
13                 nextState, reward, done, _ = self.env.step(action.clip(-2, 2))
14                 if t == self.params['MAX_STEP'] - 1:
15                     done = True
16                 ret += rewardDecay * reward
17                 rewardDecay *= self.params['gamma']
18                 stateBuf.append(state.reshape(-1))
19                 actionBuf.append(action)
20                 rewardBuf.append((reward + 8.1) / 8.1)
21
22             if steps % self.params['UPDATE_STRIDE'] == 0 or done:
23                 self.learn(nextState, done, stateBuf, actionBuf, rewardBuf)
24                 stateBuf, actionBuf, rewardBuf = [], [], []
25
26             if done:
27                 with self.totEps.get_lock():
28                     self.totEps.value += 1
29                 with self.totR.get_lock():
30                     if self.totR.value == 0:
31                         self.totR.value = ret
32                     else:
33                         self.totR.value = self.totR.value * 0.9 + ret * 0.1
34                 self.Q.put(self.totR.value)
35                 print("Rank: %d\tEps: %d\tTotRet: %f"%(self.rank,
self.totEps.value, self.totR.value))
36                 break
37
38             state = nextState
39             steps += 1
40             self.Q.put(None)
41             self.env.close()
```

So, we will interact with the environment for `MAX_EPISODE` episodes totally. In an episode (inner `for` loop), the worker interacts with the environment and record state, action and reward. Each `UPDATE_STRIDE` steps, the worker updates the global network and local network by calling `learn`. When an episode terminates, we need to update

the global counter `totEps` and record the moving average of global return value. We add lock because there is a critical section.

In the `learn` function, we calculate gradients and update networks.

```
1  def learn(self, nextState, done, stateBuf, actionBuf, rewardBuf):
2      if done:
3          R = 0
4      else:
5          R = self.LNet.forward(torch.from_numpy(nextState.reshape(1,
6          -1).astype(np.float32)).to(device))[-1][0].item()
7
8      RBuf = []
9      for r in rewardBuf[::-1]:
10         R = r + self.params['gamma'] * R
11         RBuf.append(R)
12     RBuf.reverse()
13
14     loss = self.LNet.loss(
15         torch.from_numpy(np.vstack(stateBuf).astype(np.float32)).to(device),
16         torch.from_numpy(np.vstack(actionBuf).astype(np.float32)).to(device),
17         torch.from_numpy(np.array(RBuf).reshape(-1,
18         1).astype(np.float32)).to(device)
19     )
20     self.opt.zero_grad()
21     loss.backward()
22     torch.nn.utils.clip_grad_norm_(self.LNet.parameters(), 20)
23     for l, g in zip(self.LNet.parameters(), self.GNet.parameters()):
24         g._grad = l.grad
25     self.opt.step()
26
27     self.LNet.load_state_dict(self.GNet.state_dict())
```

We first calculate the return (Q) values for each states in this period of experience. Then calculate the loss function and perform backpropagation. After it, the gradients is calculated by `torch` and we can copy them to the global network and perform update using `opt.step()`. At last, we update local network by copying parameters from the global network directly.

The final A3C implementation is shown as:

```
1  class A3C:
2      def __init__(self, gamma, updateStride, maxEps, maxSteps, hiddenSize, lr):
3          self.params = {'MAX_EPISODE': maxEps, 'MAX_STEP': maxSteps,
4          'UPDATE_STRIDE': updateStride, 'gamma': gamma, 'hiddenSize': hiddenSize}
5          self.env = gym.make('Pendulum-v0').unwrapped
6          self.globalNet = ACNet(self.env.observation_space.shape[0], hiddenSize,
7          self.env.action_space.shape[0]).to(device)
8          self.globalNet.share_memory()
9          self.opt = SharedAdam(self.globalNet.parameters(), lr=lr, betas=(0.95,
10          0.999))
11
12          self.totEps = mp.Value('i', 0)
13          self.totR = mp.Value('d', 0.0)
```

```

10     self.Q = mp.Queue()
11
12     def train(self):
13         workers = [Worker(rank, self.globalNet,
14                             ACNet(self.env.observation_space.shape[0],
15 self.params['hiddenSize'], self.env.action_space.shape[0]).to(device),
16                             self.opt,
17                             self.totEps,
18                             self.totR,
19                             self.Q,
20                             self.params)
21                             for rank in range(mp.cpu_count())]
22         [w.start() for w in workers]
23         res = []
24         while True:
25             r = self.Q.get()
26             if r is not None:
27                 res.append(r)
28             else:
29                 break
30         [w.join() for w in workers]

```

Notice that we have to make the parameters of global network shared using `globalNet.share_memory()`. In the `train()`, we need to first `start` all the workers then get some useful data using queue and at last, `join` those workers.

H2 2. Result & Discussion

After training for about 10000 episodes, we can keep the pendulum standing up once an episode starts. Wherever the pendulum is initialized, our algorithm will spin it up to the right position and apply small efforts to keep its balance. I will provide the network parameters. And you can test it using `test()` defined in the code.

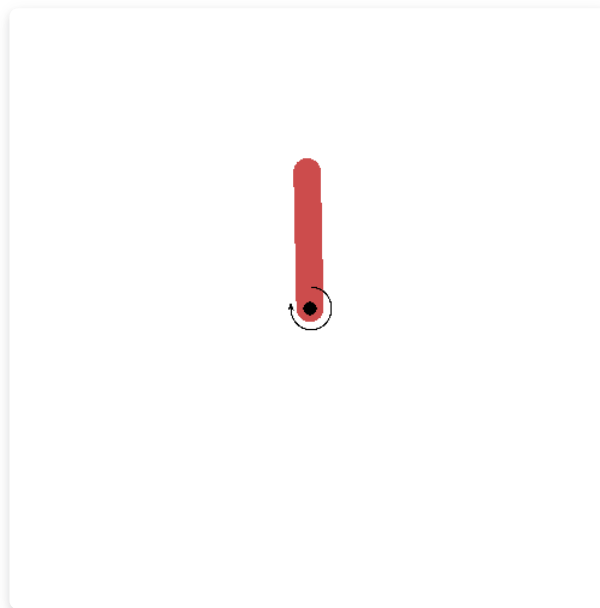


Fig.3 Result

However, it's not easy to train A3C. During my training process, I found the following problems:

1. Hyper parameter & initialization sensitivity
2. Instability
3. Cannot utilize CUDA

For the first problem, here are my final hyper parameter settings:

Hyper Parameter	Value
number of workers	8
max steps per episode	200
hidden layer size of actor network	256
hidden layer size of critic network	256
discount γ	0.99
learning rate α	0.0001
global network update stride	20
total training episodes	10000

I found that when the stride of updating global network is 5 or 10, the algorithm will hit the best episode return value of about -100 in 4000 episodes. However, after 5000 episodes, the performance dropped sharply and fell back to around -600, which is the same as it in the beginning of training. It is probably due to the short sight effect when using little experience to adjust the global network. When I set the stride to 20, it seemed the problem solved. However, it will cost more time to train, as you can see in the following figure, we first hit -100 after 7000 episodes. We also don't know if the performance will drop if we train it for more episodes. Also, if the initialization has some problem, we have to re-train it. Sometimes, the moving average return will fluctuate around -600 and the network learns a bad policy. But sometimes, we will get the result in Fig. 4.

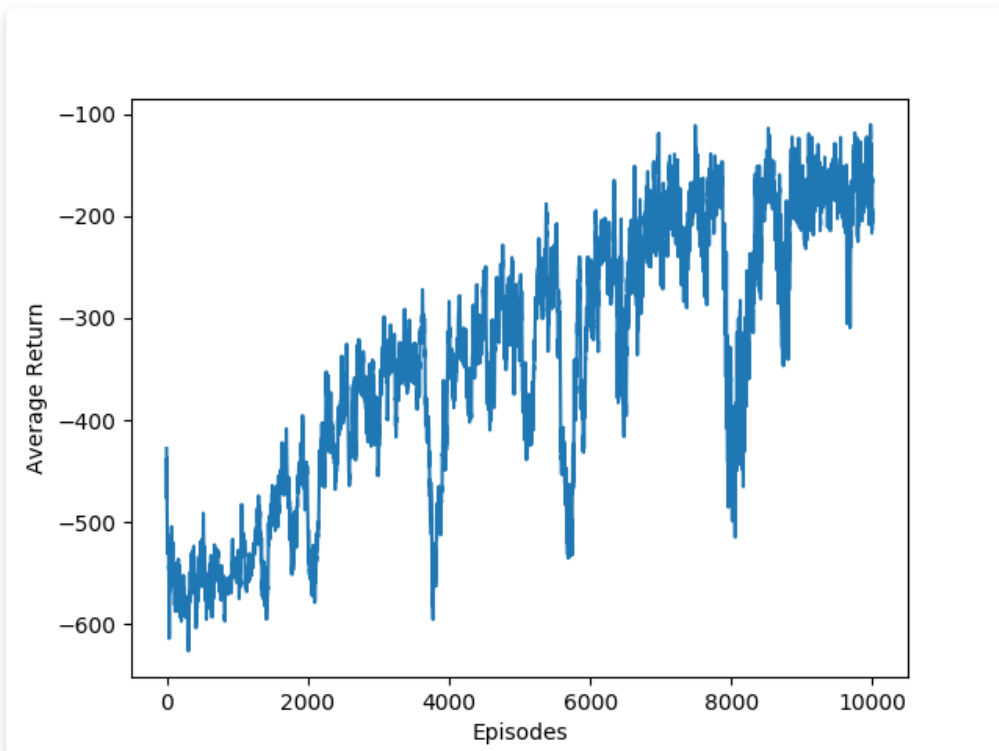


Fig.4 Moving Average of Episode Returns

As mentioned it's a bit slow to train A3C, so I reduced γ from 0.99 to 0.9 and recorded the moving average discounted return value.

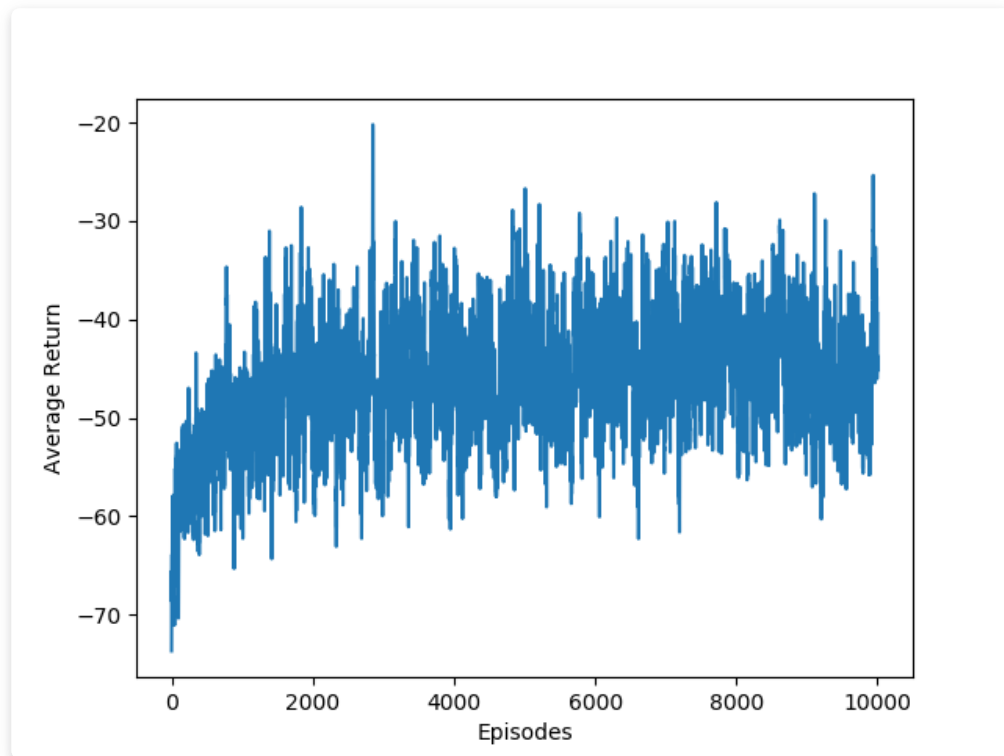


Fig.5 Moving Average of Episode Returns, gamma=0.9

I used discounted return, so the curve seems a bit messy and instable. However, you can observe that in Fig. 5, it converged in about 1500 episodes. After that point, the result fluctuates around -45 and the Pendulum can always stand up. So, the value of γ also affects the result a lot.

For the second problem, in Fig.4, you can observe some valleys in the curve. Those valleys show that the random policy always has the problem of instability especially in the circumstance that a little randomness will lead to a huge failure. If we sample the distribution, we cannot avoid sampling some outliers and make the situation worse. Fortunately, in the most of time, it will sample the right thing and thus can lead to a good result. So, I think we DDPG may perform better in Pendulum environment due to its deterministic property.

For the last problem, it's hard to run multiple workers (processes) and put all the local network to GPU. It will incur the memory allocation failure due to the memory limitation of my GPU. So, I just train those networks in CPU and, fortunately, A3C with Pendulum is not hard to run and I got good results.

H3 3. Acknowledgement

When implementing my A3C algorithm, I referred to the following GitHub repositories, which helped me a lot. I learned the network structure, the loss function and also coding techniques using `torch` from those codes.

- <https://github.com/MorvanZhou/pytorch-A3C>
- <https://github.com/ikostrikov/pytorch-a3c>