

1 Introduction

- 1.1 Actor-critic (AC) Algorithm
- 1.2 Asynchronous Advantage Actor-critic (A3C) Algorithm
- 1.3 Deep Deterministic Policy Gradient (DDPG) Algorithm

2 Environment configuration: *Pendulum-v1*

3 Experimental contents

- 3.1 Code structure for A3C
- 3.2 Network structure and parameter settings for A3C
- 3.3 Code structure for DDPG
- 3.4 Network structure and parameter settings for DDPG

4 Analysis of Results

- 4.1 The result of A3C
 - 4.1.1 Maximum number of steps within an *Episode MAX_EP_STEP*
 - 4.1.2 Number of *Worker* threads
- 4.2 The result of DDPG

5 Experimental Insights

1 Introduction

Algorithms such as deep Q networks have achieved good performance in solving complex tasks with unprocessed, high-dimensional inputs. However, *DQN* is only able to solve problems that are discrete and located in low-dimensional action spaces due to the need to find discrete actions that maximize *Q*. With the continuity property of *policy gradient*, the *Actor-Critic* algorithm can be used to solve problems with continuous and high-dimensional action spaces. In this paper, we present two improved algorithms based on the AC algorithm: the **Asynchronous Advantage Actor-critic (A3C)** algorithm and the **Deep Deterministic Policy Gradient (DDPG) Algorithm**, with experiments and performance analysis in the *gym* environment of *Pendulum-v0* (*Pendulum-v1*). Since in our implemented version of *gym*, *Pendulum-v0* has already been discarded and replaced by *Pendulum-v1* (actually the same thing), so in the next following sections *Pendulum-v0* and *Pendulum-v1* means the same thing.

1.1 Actor-critic (AC) Algorithm

Actor-Critic typically consists of two parts, the actor (*Actor*) and the evaluator (*Critic*). For the first part, *Actor* uses policy functions and is responsible for generating actions (*Action*) and interacting with the environment. And for the second part, *Critic* uses the value function, which is responsible for evaluating the performance of *Actor* and guiding the next stage of the Actor's action. In the *Actor-Critic* algorithm, we need to make two sets of approximations, which have also been mentioned in *Policy-based Reinforcement Learning*.

- The first set is the approximation of the policy function.

$$\pi_{\theta}(s, a) = P(a|s, \theta) \approx \pi(a|s)$$

- The approximation of the value function at the second group, for the state value and action value functions, respectively, is

$$\hat{v}(s, w) \approx v_{\pi}(s)$$

$$\hat{q}(s, a, w) \approx q_{\pi}(s, a)$$

In the previously learned *Monte-Carlo Policy Gradient (REINFORCE)* algorithm, we update the network parameters by the stochastic gradient ascent *Stochastic Gradient Ascent* method, referring to the policy gradient theorem, and use the return value v_t to as the *unbiased sample* of $Q^{\pi_{\theta}}(s_t, a_t)$, yielding the parameter update formula of the strategy as

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) v_t$$

For *Actor* in the AC system, the above formula for policy update is needed to keep updating the parameters and generating optimal actions. However, unlike *Monte-Carlo Policy Gradient (REINFORCE)*, the v_t in the formula no longer comes from Monte-Carlo sampling, but should be obtained from another component of the AC system, which is *Critic*. For *Critic* part, we refer to *DQN* and use a Q network for the representation, where the input can be the state and the output is the value of the optimal action or the respective value of all actions.

Thus, the main framework of the AC algorithm is.

(1) *Critic* computes the optimal value v_t of the current state through the Q network and returns it to *Actor*; after getting the feedback and the new state from *Actor*, it updates the Q network

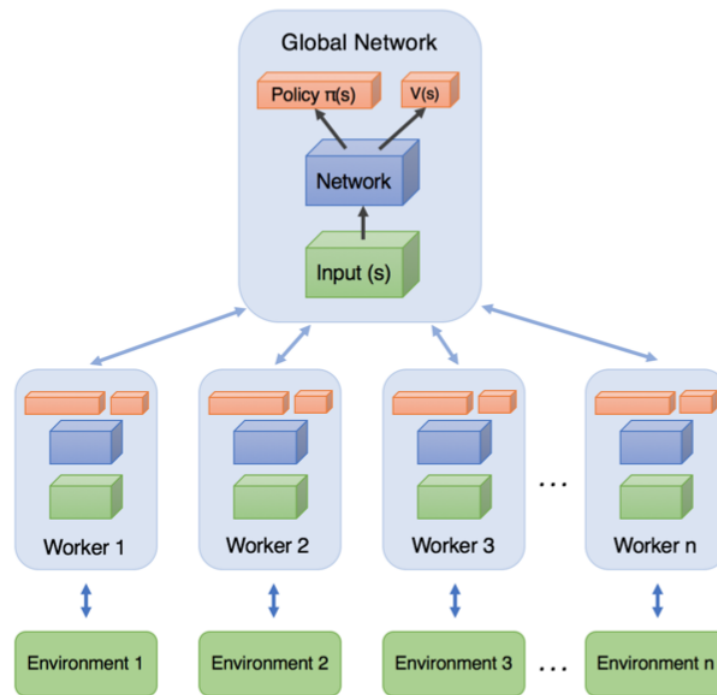
(2) *Actor* iteratively updates the argument θ of the policy function using v_t and selects the action with it, gets the feedback and the new state, and returns it to *Critic*

1.2 Asynchronous Advantage Actor-critic (A3C) Algorithm

In the use of the AC algorithm, although single-step updates are possible, which are faster than the round updates of the traditional *Policy Gradient* algorithm, the algorithm itself is difficult to converge because the behavior of *Actor* depends on the *Value* of *Critic*, and since *Critic* is inherently difficult to converge. The reason is that the behavior of *Actor* depends on the *Value* of *Critic*.

To address the above problem, we can use empirical replay to solve the hard-to-converge problem by referring to the implementation of the *DQN* algorithm. However, there is a problem with the experience replay itself, as the experience data in the replay pool is too correlated and is likely to be ineffective when used for training.

The A3C algorithm solves the problem of experience replay data correlation based on the idea of experience replay, using a multi-threaded method of learning by interacting with the environment simultaneously. In the multi-thread, each thread aggregates the successes learned from its own interaction with the environment, organizes and saves them in a common place, and periodically brings out the results of everyone's concerted learning to guide its own learning interactions with the environment later, in order to achieve the goal of **asynchronous parallelism**. The network structure of the A3C algorithm is shown in the figure below.



The A3C algorithm was first proposed by *DeepMind* to solve the problem of difficult convergence of the *Actor-Critic* algorithm. The A3C algorithm exploits the power of parallel computing by creating multiple parallel environments in which multiple *workers* with the same AC network structure (*local network*) interact independently, store their experiences, and perform periodic parameter updates to the *Global AC Network* simultaneously. While the parameters are updated, each *worker* also updates its own *local network* with the network parameters of the *global network*, thus reducing the correlation of parameter updates and thus improving the convergence of the whole algorithm on the basis of guaranteed iterative learning. The pseudo-code of the A3C algorithm is as follows.

Algorithm 1: Asynchronous advantage actor-critic for each actor-learner thread

```
1 Assume global shared parameter vectors  $\theta$  and  $\theta_v$  and global shared counter  $T = 0$ ;  
2 Assume thread-specific parameter vector  $\theta'$  and  $\theta'_v$ ;  
3 Initialize thread step counter  $t \leftarrow 1$ ;  
4 while true do  
5   Reset gradients:  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ ;  
6   Synchronize thread-specific parameters  $\theta' = \theta$  and  $\theta'_v = \theta_v$ ;  
7    $t_{start} = t$  ;  
8   Get state  $s_t$ ;  
9   while true do  
10    Perform  $a_t$  according to policy  $\pi(a_t|s_t; \theta')$ ;  
11    Receive reward  $r_t$  and new state  $s_{t+1}$ ;  
12     $t \leftarrow t + 1$ ;  
13     $T \leftarrow T + 1$ ;  
14    if terminal  $s_t$  or  $t - t_{start} == t_{max}$  then  
15      Break;  
16    end  
17  end  
18   $R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{ Bootstrap from last state} \end{cases}$  ;  
19  for  $i \in \{t - 1, \dots, t_{start}\}$  do  
20     $R \leftarrow r_i + \gamma R$ ;  
21    Accumulate gradients wrt  $\theta'$ :  $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta') (R - V(s_i; \theta'_v))$ ;  
22    Accumulate gradients wrt  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + \partial (R - V(s_i; \theta'_v))^2 / \partial \theta'_v$ ;  
23  end  
24  Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ .;  
25  if  $T > T_{max}$  then  
26    Break;  
27  end  
28 end
```

Compared to the traditional AC, the A3C algorithm has three main improvements.

- **Asynchronous training framework.**

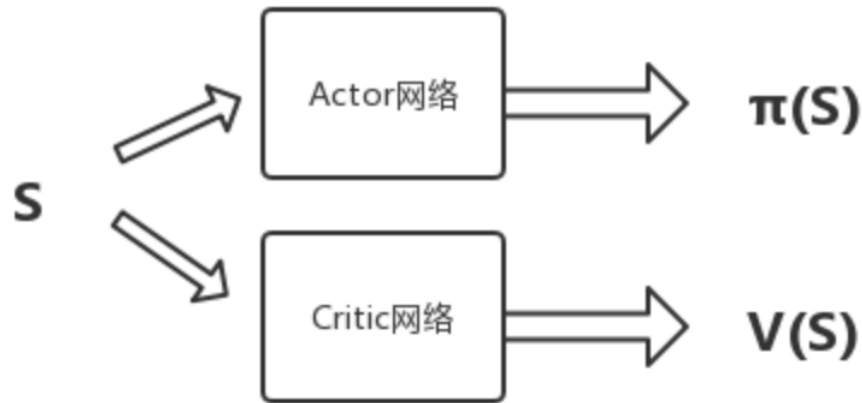
The **Global Network** above in the previous figure is the shared common part stated in 1.2 section, which is mainly a common neural network model. This neural network** includes the functions of both *Actor* network and *Critic* network. **There are n worker threads**, each of which has the same network structure as the common neural network, **and each thread will independently interact with the environment to get the experience data, and these threads will not interfere with each other and run independently****.

After interacting with the environment for a certain amount of data, each thread computes the gradient of the loss function of the neural network in its own thread, but these gradients do not update the neural network in its own thread, but update the public neural network. That is, **n threads independently update the parameters of the common part of the neural network model** using the accumulated gradients. ** Every once in a while, the threads update the parameters of their own neural network to those of the public neural network, which in turn guides the later environmental interactions**.

As can be seen, the public part of the network model is the model we want to learn, while the network model in the thread is mainly used for interaction with the environment. These models in the thread can help the thread to better interact with the environment and get high-quality data to help the model converge faster.

- **Network Structure Optimization.**

Unlike the AC algorithm, here we logically merge the *Actor* and *Critic* networks, i.e., input state, output state value v_t and the corresponding policy π , but of course, the binaries are still physically two separate networks, handled separately as follows



- Optimization of **Critic evaluation points**.

In the AC algorithm, we discuss the use of single-step sampling to approximate the estimation of $Q(S, A)$, i.e.: $Q(S, A) = R + \gamma V(S')$, where the value of $V(S)$ needs to be obtained by *Critic* network learning. Thus the dominance function can be expressed as

$$A(S, t) = R + \gamma V(S') - V(S)$$

In the A3C algorithm, **N* step sampling is used to speed up convergence**, so the dominance function can be expressed as

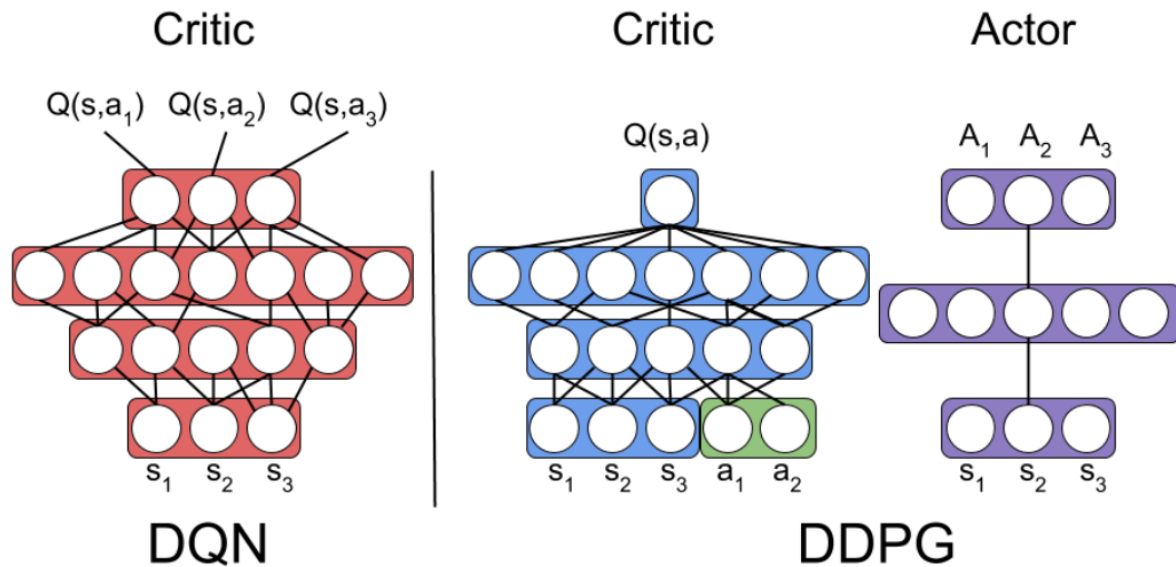
$$A(S, t) = R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n V(S') - V(S)$$

The loss function part for *Actor* and *Critic* is essentially the same as for AC. One small optimization point is to add the entropy term of the strategy π to the loss function of the *Actor-Critic* strategy function with a coefficient of c , i.e., the gradient update of the strategy parameters becomes the following form.

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) A(S, t) + c \nabla_{\theta} H(\pi(S_t, \theta))$$

1.3 Deep Deterministic Policy Gradient (DDPG) Algorithm

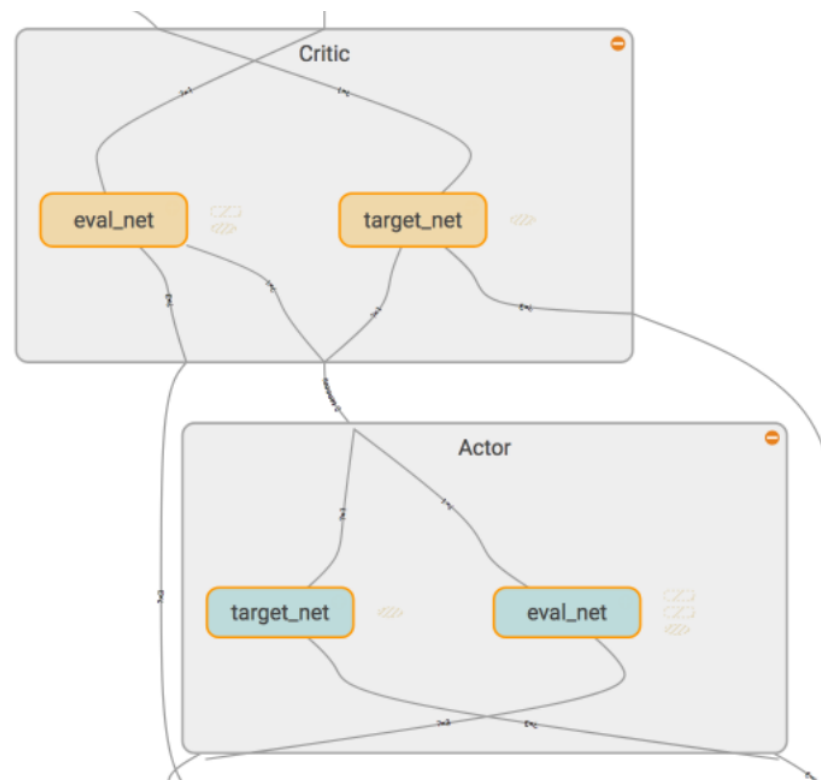
Deep Deterministic Policy Gradient (DDPG) is a model-free off-policy actor-critic algorithm that combines *DQN* and *DPG*. As it is an off-policy algorithm, *DDPG* uses two separate policies for the exploration and updates. It uses a stochastic behaviour policy for the exploration, and deterministic policy for the target update.



- Network Schematics

As another improved algorithm based on the AC algorithm, *DDPG* has two networks: actor and critic. The actor produces the action to explore, and during the updating of the actor, TD error from a critic is calculated, and the critic network is updated based on the it which is similar to Q-learning update rule.

In *DDPG*, four neural networks are used: a Q network, a deterministic policy network, a target Q network, and a target policy network. For the Q network and policy network, the actor directly maps the states to actions, instead of outputting the probability distribution across a discrete action space, which is an important feature of *DDPG*. For the target networks, they are time-delayed copies of their original networks, which slowly track the learned networks. With the target networks, the stability in learning can be safeguarded.



- Learning

The pseudo-code of the algorithm is below:

Algorithm 1 Deep Deterministic Policy Gradient

- 1: Input: initial policy parameters θ , Q-function parameters ϕ , empty replay buffer \mathcal{D}
- 2: Set target parameters equal to main parameters $\theta_{\text{targ}} \leftarrow \theta$, $\phi_{\text{targ}} \leftarrow \phi$
- 3: **repeat**
- 4: Observe state s and select action $a = \text{clip}(\mu_\theta(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$, where $\epsilon \sim \mathcal{N}$
- 5: Execute a in the environment
- 6: Observe next state s' , reward r , and done signal d to indicate whether s' is terminal
- 7: Store (s, a, r, s', d) in replay buffer \mathcal{D}
- 8: If s' is terminal, reset environment state.
- 9: **if** it's time to update **then**
- 10: **for** however many updates **do**
- 11: Randomly sample a batch of transitions, $B = \{(s, a, r, s', d)\}$ from \mathcal{D}
- 12: Compute targets

$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

- 13: Update Q-function by one step of gradient descent using

$$\nabla_\phi \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_\phi(s, a) - y(r, s', d))^2$$

- 14: Update policy by one step of gradient ascent using

$$\nabla_\theta \frac{1}{|B|} \sum_{s \in B} Q_\phi(s, \mu_\theta(s))$$

- 15: Update target networks with

$$\begin{aligned}\phi_{\text{targ}} &\leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta\end{aligned}$$

- 16: **end for**
 - 17: **end if**
 - 18: **until** convergence
-

As in many other algorithms, *DDPG* also uses a replay buffer to sample experience to update neural network parameters. During each trajectory roll-out, we save all the experience tuples (state, action, reward, next_state) and store them in a finite-sized cache — the “replay buffer.” Then, we sample random mini-batches of experience from the replay buffer when we update the value and policy networks.

The value network is updated similarly as is done in Q-learning. As we calculate the next-state Q values with the target value network and target policy network, we minimize the mean-squared loss between the updated Q value and the original Q value, which is generally known as TD error:

$$Loss = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$$

For the policy function, our objective is to maximize the expected return:

$$J(\theta) = \mathbb{E}[Q(s, a) |_{s=s_t, a=\mu(s_t)}]$$

Since we are updating the policy in an off-policy way with batches of experience, we take the mean of the sum of gradients calculated from the mini-batch.

For the target networks, we have them slowly track those of the learned networks via “soft updates” as illustrated below:

$$\theta^{Q'} = \tau\theta^Q + (1 - \tau)\theta^{Q'}$$

$$\theta^{\mu'} = \tau\theta^{\mu} + (1 - \tau)\theta^{\mu'}$$

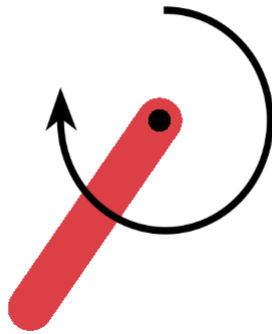
- Exploration

In Reinforcement learning for discrete action spaces, exploration is done via probabilistically selecting a random action. For continuous action spaces, exploration is done via adding noise to the action itself.

$$\mu'(s_t) = \mu(s_t|\theta_t^{\mu}) + \mathcal{N}(0, \sigma)$$

2 Environment configuration: *Pendulum-v1*

The scenario for the *Pendulum-v1* problem has a frictionless, vertically swingable pendulum, as shown in the image below on the left. The goal of the problem is to train *agent* so that the pendulum can point as far up as possible and remain balanced. The pendulum starts at a random position each time, and since the scenario is an *unsolved environment*, there is no reward threshold and no explicit termination condition such as "end at a certain point". The scenario for the *Pendulum-v1* problem has three observed inputs: the sine of the pendulum angle $\sin(\theta)$, the cosine $\cos(\theta)$, and the angular velocity θ_{dt} , with the three inputs taking values in the range shown on the right below.



Observation	Min	Max
$\sin(\theta)$	-1.0	+1.0
$\cos(\theta)$	-1.0	+1.0
θ_{dt}	-8.0	+8.0

The joint left/right effect of the three inputs takes values in the range $[-2, 2]$ and the estimation function for *reward* is as follows

$$R = -(\theta^2 + 0.1 \times \theta_{dt}^2 + 0.001 \times action^2)$$

where θ takes values in the range of $[-\pi, \pi]$ and θ_{dt} represents the angular velocity of the pendulum. For this experiment, we use *gym* to simulate this environment, and a simple environment configuration is shown below.

```

1 import gym
2 env = gym.make('Pendulum-v1')
3 env.reset()
4 for _ in range(1000):
5     env.render()
6     env.step(env.action_space.sample()) # take a random action
7 env.close()
```


3 Experimental contents

3.1 Code structure for A3C

The implementation of this algorithm is divided into two classes and the main function. The two classes namely as *Class ACNet* and *Class Worker*.

The *ACNet* class mainly builds the basic architecture of the AC network, including network design, initialization, action selection and loss function calculation.

The *ACNet* class first make `__init__` part, which is shown as follows:

```
1  def __init__(self, scope, globalAC=None):
2
3      if scope == GLOBAL_NET_SCOPE:  # if global network
4          with tf.variable_scope(scope):
5              self.s = tf.placeholder(tf.float32, [None, N_S], 's')
6              self.a_params, self.c_params = self._build_net(scope)[-2:]
7      else:  # else local net
8          with tf.variable_scope(scope):
9              self.s = tf.placeholder(tf.float32, [None, N_S], 's')
10             self.a_his = tf.placeholder(tf.float32, [None, N_A], 'A')
11             self.v_target = tf.placeholder(tf.float32, [None, 1], 'Vtarget')
12
13             mu, sigma, self.v, self.a_params, self.c_params =
self._build_net(scope)
14
15             with tf.name_scope('wrap_a_out'):
16                 mu = mu * A_BOUND[1]
17                 sigma = sigma + 1e-4
18
19             normal_dist = tf.distributions.Normal(mu, sigma)
20
21             td = tf.subtract(self.v_target, self.v, name='TD_error')
22
23             with tf.name_scope('a_loss'):
24                 log_prob = normal_dist.log_prob(self.a_his)
25                 exp_v = log_prob * tf.stop_gradient(td)
26                 entropy = normal_dist.entropy()
27                 self.exp_v = ENTROPY_FACTOR * entropy + exp_v
28                 self.a_loss = tf.reduce_mean(-self.exp_v)
29
30             with tf.name_scope('c_loss'):
31                 self.c_loss = tf.reduce_mean(tf.square(td))
32
33             with tf.name_scope('local_grad'):
34                 self.a_grads = tf.gradients(self.a_loss, self.a_params)
35                 self.c_grads = tf.gradients(self.c_loss, self.c_params)
36
37             with tf.name_scope('choose_a'):  # use local params to choose
action
38                 self.A = tf.clip_by_value(tf.squeeze(normal_dist.sample(1),
axis=[0, 1]), A_BOUND[0], A_BOUND[1])
39
40             with tf.name_scope('sync'):
41                 with tf.name_scope('pull'):
```

```

42         self.pull_a_params_op = [l_p.assign(g_p) for l_p, g_p in
zip(self.a_params, globalAC.a_params)]
43         self.pull_c_params_op = [l_p.assign(g_p) for l_p, g_p in
zip(self.c_params, globalAC.c_params)]
44         with tf.name_scope('push'):
45             self.update_a_op =
OPT_ACTOR.apply_gradients(zip(self.a_grads, globalAC.a_params))
46             self.update_c_op =
OPT_CRITIC.apply_gradients(zip(self.c_grads, globalAC.c_params))

```

Divided into two parts, the first part is for both global net and local net, used to set variables and call another function `_build_net`. Loss function is also defined with following code:

```

1  with tf.name_scope('a_loss'):
2      log_prob = normal_dist.log_prob(self.a_hist)
3      exp_v = log_prob * tf.stop_gradient(td)
4      entropy = normal_dist.entropy()
5      self.exp_v = ENTROPY_FACTOR * entropy + exp_v
6      self.a_loss = tf.reduce_mean(-self.exp_v)

```

```

1  with tf.name_scope('c_loss'):
2      self.c_loss = tf.reduce_mean(tf.square(td))

```

We also defined `local_grad` for gradient update.

```

1  with tf.name_scope('local_grad'):
2      self.a_grads = tf.gradients(self.a_loss, self.a_params)
3      self.c_grads = tf.gradients(self.c_loss, self.c_params)

```

We also use local params to choose action. Based on states, we generate probability distributions for actions and sample.

```

1  with tf.name_scope('choose_a'): # use local params to choose action
2      self.A = tf.clip_by_value(tf.squeeze(normal_dist.sample(1), axis=[0, 1]),
A_BOUND[0], A_BOUND[1])

```

We also do updates for global network and local network with `pull` and `push` operations.

```

1  with tf.name_scope('sync'):
2      with tf.name_scope('pull'):
3          self.pull_a_params_op = [l_p.assign(g_p) for l_p, g_p in
zip(self.a_params, globalAC.a_params)]
4          self.pull_c_params_op = [l_p.assign(g_p) for l_p, g_p in
zip(self.c_params, globalAC.c_params)]
5          with tf.name_scope('push'):
6              self.update_a_op = OPT_ACTOR.apply_gradients(zip(self.a_grads,
globalAC.a_params))
7              self.update_c_op = OPT_CRITIC.apply_gradients(zip(self.c_grads,
globalAC.c_params))

```

The second part for ACNet is `_build_net`. Its implementation is as follows:

```

1  def _build_net(self, scope):
2      w_init = tf.random_normal_initializer(0., .1)
3      with tf.variable_scope('actor'):
4          l_a = tf.layers.dense(self.s, 200, tf.nn.relu6,
kernel_initializer=w_init, name='l_a')
5          mu = tf.layers.dense(l_a, N_A, tf.nn.tanh,
kernel_initializer=w_init, name='mu')
6          sigma = tf.layers.dense(l_a, N_A, tf.nn.softplus,
kernel_initializer=w_init, name='sigma')
7          with tf.variable_scope('critic'):
8              l_c = tf.layers.dense(self.s, 100, tf.nn.relu6,
kernel_initializer=w_init, name='l_c')
9              v = tf.layers.dense(l_c, 1, kernel_initializer=w_init, name='v') #
state value
10
11         a_params = tf.get_collection(tf.GraphKeys.TRAINABLE_VARIABLES,
scope=scope + '/actor')
12         c_params = tf.get_collection(tf.GraphKeys.TRAINABLE_VARIABLES,
scope=scope + '/critic')
13         return mu, sigma, v, a_params, c_params

```

This part includes both actor network and critic network. After that we defined some functions necessary in following code:

```

1  def choose_action(self, s): # run by a local
2      s = s[np.newaxis, :]
3      return SESS.run(self.A, {self.s: s})
4
5  def pull_global(self): # run by a local
6      SESS.run([self.pull_a_params_op, self.pull_c_params_op])
7
8  def update_global(self, input_dict): # run by a local
9      SESS.run([self.update_a_op, self.update_c_op], input_dict) # apply local
grads to global net

```

The *Worker* class creates a separate *Pendulum-v1* environment for each *worker thread*, calls the *ACNet* class to build a *local network* with `__init__`

`__init__()`: declares variables

```

1  def __init__(self, name, globalAC):
2      self.env = gym.make(GAME).unwrapped
3      self.name = name
4      self.AC = ACNet(name, globalAC)

```

Then `work` sets up the environment, training episodes, and records rewards during training for further analysis.

```

1  def work(self):
2      global GLOBAL_RUNNING_R, GLOBAL_EP
3      buffer_s, buffer_a, buffer_r = [], [], []
4      total_step = 1
5
6      while not COORDINATOR.should_stop() and GLOBAL_EP < MAX_GLOBAL_EP:

```

```

7         s = self.env.reset()
8         ep_r = 0
9         for ep_t in range(MAX_EP_STEP):
10             a = self.AC.choose_action(s)
11             s_, r, done, info = self.env.step(a)
12
13             ep_r += r
14             buffer_s.append(s)
15             buffer_a.append(a)
16             buffer_r.append((r+8)/8)    # normalize the reward
17
18             if ep_t == MAX_EP_STEP - 1:
19                 done = True
20             else:
21                 done = False
22
23             if total_step % UPDATE_GLOBAL_ITER == 0 or done:    # update
global and assign to local net
24                 if done:
25                     v_s_ = 0
26                 else:
27                     v_s_ = SESS.run(self.AC.v, {self.AC.s: s_[np.newaxis,
:]})[0, 0]
28
29                 buffer_v_target = []
30                 for r in buffer_r[::-1]:    # reverse buffer r
31                     v_s_ = r + GAMMA * v_s_
32                     buffer_v_target.append(v_s_)
33
34                 buffer_v_target.reverse()
35
36                 buffer_s, buffer_a, buffer_v_target = np.vstack(buffer_s),
np.vstack(buffer_a), np.vstack(buffer_v_target)
37                 input_dict = {
38                     self.AC.s: buffer_s,
39                     self.AC.a_hist: buffer_a,
40                     self.AC.v_target: buffer_v_target,
41                 }
42                 self.AC.update_global(input_dict)
43                 buffer_s, buffer_a, buffer_r = [], [], []
44                 self.AC.pull_global()
45
46                 s = s_
47                 total_step += 1
48                 if done:
49                     if len(GLOBAL_RUNNING_R) == 0:    # record running episode
reward
50                         GLOBAL_RUNNING_R.append(ep_r)
51                     else:
52                         GLOBAL_RUNNING_R.append(0.9 * GLOBAL_RUNNING_R[-1] + 0.1
* ep_r)
53
54                     print(
55                         self.name,
56                         "Ep:", GLOBAL_EP,
57                         "| Ep_reward: %i" % GLOBAL_RUNNING_R[-1],
58                         )
59                     GLOBAL_EP += 1
60                     break

```

The third part for the code is `main` part, we set up the training device and the training optimizer to be *RMSProp*. We also set global net and use `work` to join the pre-set number of worker threads for training.

```
1  if __name__ == "__main__":
2      SESS = tf.Session(config=tf.ConfigProto(log_device_placement=True))
3      COORDINATOR = tf.train.Coordinator()
4
5      with tf.device("/gpu:0"):
6          OPT_ACTOR = tf.train.RMSPropOptimizer(LR_ACTOR, name='RMSPropA')
7          OPT_CRITIC = tf.train.RMSPropOptimizer(LR_CRITIC, name='RMSPropC')
8          GLOBAL_AC = ACNet(GLOBAL_NET_SCOPE)
9
10         workers = []
11         for i in range(N_WORKERS):
12             i_name = 'w_%i' % i
13             workers.append(worker(i_name, GLOBAL_AC))
14
15         SESS.run(tf.global_variables_initializer())
16
17         if OUTPUT_GRAPH:
18             if os.path.exists(LOG_DIR):
19                 shutil.rmtree(LOG_DIR)
20             tf.summary.FileWriter(LOG_DIR, SESS.graph)
21
22         worker_threads = []
23         for worker in workers:
24             job = lambda: worker.work()
25             t = threading.Thread(target=job)
26             t.start()
27             worker_threads.append(t)
28
29         COORDINATOR.join(worker_threads)
```

At last we make use of matplotlib to draw figures for our training and reward changing process.

```
1  plt.plot(np.arange(len(GLOBAL_RUNNING_R)), GLOBAL_RUNNING_R)
2  plt.xlabel('Step')
3  plt.ylabel('Total Reward')
4  plt.title('A3C results')
5  plt.show()
```

3.2 Network structure and parameter settings for A3C

In the neural network implemented by the A3C algorithm, since the states in the *Pendulum-v1* environment are abstracted as $\sin(\theta)$, $\cos(\theta)$ and angular velocity θ_{dt} of the pendulum with low feature dimension, there is no need to use convolutional layers to implement the hidden layer, but A simple design with several fully connected layers is still necessary for implementing the action selection network for *Actor* and the state value computation network for *Critic*, respectively.

The input of this sub-network is the state *state*, and the output is a distribution of action values (*mu*, *sigma*) that characterizes the prediction of *Actor* in the current state. In the subsequent *choose action* function, a normal distribution is constructed based on these two statistics. Normal distribution from *action* is sampled. The purpose of using such an approach for action selection is to better accommodate action continuity in the *Pendulum-v1* environment. After testing, the activation functions of the fully connected layers are set to *relu6*, *tanh* and *softplus*, noting that in order to make the variance of the normal distribution larger than 0, we need to add a *bias* = $1e-4$ in the calculation of *sigma*.

```
1 with tf.variable_scope('actor'):
2     l_a = tf.layers.dense(self.s, 200, tf.nn.relu6,
3         kernel_initializer=w_init, name='l_a')
4     mu = tf.layers.dense(l_a, N_A, tf.nn.tanh, kernel_initializer=w_init,
5         name='mu')
6     sigma = tf.layers.dense(l_a, N_A, tf.nn.softplus,
7         kernel_initializer=w_init, name='sigma')
```

```
1 with tf.name_scope('wrap_a_out'):
2     mu = mu * A_BOUND[1]
3     sigma = sigma + 1e-4
```

Critic network: the fully connected layer is set with 100 neurons. The input of this sub-network is the state *state* and the output is the state value of that state evaluated by *Critic*. The activation function of the fully connected layer is set to *relu6*.

```
1 with tf.variable_scope('critic'):
2     l_c = tf.layers.dense(self.s, 100, tf.nn.relu6,
3         kernel_initializer=w_init, name='l_c')
4     v = tf.layers.dense(l_c, 1, kernel_initializer=w_init, name='v') # state
5     value
```

The entire model is set to have a *discount factor*- *GAMMA* of 0.9, an *entropy factor*- *ENTROPY_FACTOR* of 0.01, a *global network* that is updated every 10 rounds with parameter *UPDATE_GLOBAL_ITER*, a default maximum step count *MAX_EP_EP* of 200, a total *episode* count *MAX_GLOBAL_EP* of 2000. The default thread count is set to be *cpu_count*. And the learning rate for actor and critic is set to be 0.0001 and 0.001 respectively. Note that a comparative experiment on the effect of *MAX_EP_STEP* and the number of threads *N_WORKERS* on the training results will be done in the subsequent analysis of the results in Section 4 for a more comprehensive discussion of the factors affecting the experimental effect.

```
1 MAX_EP_STEP = 200 # Total steps per episode
2 MAX_GLOBAL_EP = 2000 # Total number of episodes
3 UPDATE_GLOBAL_ITER = 10
4 GAMMA = 0.9
5 LR_ACTOR = 0.0001
6 LR_CRITIC = 0.001
7 ENTROPY_FACTOR = 0.01
8 GLOBAL_EP = 0
```

3.3 Code structure for DDPG

The main class of the DDPG algorithm is the class `DDPG`. It first make `init` part, which is shown as follows. In this part, the four networks, the replay buffer and the optimizers of both actor and critic are initialized.

```
1 def __init__(self):
2     self.actor_eval_net, self.actor_target_net = Net_A(), Net_A()
3     self.critic_eval_net, self.critic_target_net = Net_C(), Net_C()
4
5     self.learn_step_counter = 0
6     self.memory_counter = 0
7     self.memory = np.zeros((MEMORY_CAPACITY, N_STATES * 2 + N_ACTIONS + 1))
8     self.actor_optimizer = torch.optim.Adam(
9         self.actor_eval_net.parameters(), lr=LR_A)
10    self.critic_optimizer = torch.optim.Adam(
11        self.critic_eval_net.parameters(), lr=LR_C)
12
13    self.loss_func = nn.MSELoss()
```

The `choose_action()`, `store_transition()` and `update_params()` functions are shown as follows.

```
1 def choose_action(self, s):
2     s = torch.unsqueeze(torch.FloatTensor(s), 0)
3     action = self.actor_eval_net(s)[0].detach()
4     return action
5
6 def store_transition(self, s, a, r, s_):
7     transition = np.hstack((s, [a, r], s_))
8     index = self.memory_counter % MEMORY_CAPACITY
9     self.memory[index, :] = transition
10    self.memory_counter += 1
11
12 def update_params(self, t, s, isSoft=False):
13     for t_param, param in zip(t.parameters(), s.parameters()):
14         if isSoft:
15             t_param.data.copy_(
16                 t_param.data * (1.0 - TAU) + param.data * TAU)
17         else:
18             t_param.data.copy_(param.data)
```

The `learn()` function is the main part of the *DDPG* algorithm, which is shown as follows. First of all, a mini-batch data is sampled from buffer, and the (s, a, r, s_-) is extracted from the sampled mini-batch.

```
1 def learn(self):
2
3     sample_index = np.random.choice(MEMORY_CAPACITY, BATCH_SIZE)
4     b_memory = self.memory[sample_index, :]
5     b_s = torch.FloatTensor(b_memory[:, :N_STATES])
6     b_a = torch.FloatTensor(b_memory[:, N_STATES:N_STATES+N_ACTIONS])
7     b_r = torch.FloatTensor(b_memory[:, -N_STATES-1:-N_STATES])
8     b_s_ = torch.FloatTensor(b_memory[:, -N_STATES:])
```

After that, an action is made and its action values is evaluated. The loss of actor network is optimized according to the equation above in section 1.3. For the critic, the target Q value is computed using the information of next state, and the TD error, which is the loss of the critic network, is calculated and optimized according to the equation above in section 1.3.

```

1     a = self.actor_eval_net(b_s)
2     q = self.critic_eval_net(b_s, a)
3     a_loss = -torch.mean(q)
4
5     self.actor_optimizer.zero_grad()
6     a_loss.backward()
7     self.actor_optimizer.step()
8
9     a_target = self.actor_target_net(b_s_)
10    q_tmp = self.critic_target_net(b_s_, a_target)
11    q_target = b_r + GAMMA * q_tmp
12    q_eval = self.critic_eval_net(b_s, b_a)
13    td_error = self.loss_func(q_target, q_eval)
14
15    self.critic_optimizer.zero_grad()
16    td_error.backward()
17    self.critic_optimizer.step()

```

Finally, the target networks are updated softly using the tool function `update_params()`, according to the equation above in section 1.3.

```

1     self.update_params(self.actor_target_net, self.actor_eval_net, True)
2     self.update_params(self.critic_target_net, self.critic_eval_net, True)

```

For the training part, in each episode, the agent is restricted to run at most `max_step_num` steps. For each step, first of all an action is chosen by the `choose_action()` function. After that, a gradually decreased Gaussian noise is added to guarantee the exploration. The $(s, a, r, s_)$ is stored by `store_transition()`. Every time when the replay buffer is full do the `learn()`. When each episode is end, store the return list.

```

1     def train(episode_num, max_step_num):
2         ddpq = DDPG()
3         return_list = []
4         action_low = env.action_space.low
5         action_high = env.action_space.high
6         var = 3.0
7         for i in range(episode_num):
8             s = env.reset()
9             ep_r = 0
10            step = 0
11            while True:
12                step += 1
13                # env.render()
14                a = ddpq.choose_action(s)
15                a = np.clip(np.random.normal(a, var), action_low, action_high)
16                s_, r, done, info = env.step(a)
17                ddpq.store_transition(s, a, (r+8)/8, s_)
18
19
20            ep_r += r

```



```

21         if ddpq.memory_counter > MEMORY_CAPACITY:
22             var *= GAMMA_VAR
23             ddpq.learn()
24         if done or step == max_step_num - 1:
25             if len(return_list) == 0:
26                 return_list.append(ep_r)
27             else:
28                 return_list.append(0.95 * return_list[-1] + 0.05 * ep_r)
29             print('Episode: ', i, ' Reward: %i' % (return_list[-1]))
30             break
31         s = s_
32     np.save('return_list_'+ str(max_step_num), return_list)

```

3.4 Network structure and parameter settings for DDPG

The actor and critic networks are shown below. For the actor network, the input is the state s and the output is an action in range $[-2, 2]$. For the critic network, the input is the state s and the action a , and the output is the evaluated $Q(s, a)$. The specific structure of the network is relatively simple and can be viewed very clearly in the code, so there is no further elaboration.

```

1  class Net_A(nn.Module):
2      def __init__(self, ):
3          super(Net_A, self).__init__()
4          self.fc1 = nn.Linear(N_STATES, 50)
5          self.fc1.weight.data.normal_(0, 0.1)
6          self.out = nn.Linear(50, N_ACTIONS)
7          self.out.weight.data.normal_(0, 0.1)
8
9      def forward(self, s):
10         x = self.fc1(s)
11         x = F.relu(x)
12         x = self.out(x)
13         x = torch.tanh(x)
14         action = x * 2
15         return action
16
17
18  class Net_C(nn.Module):
19      def __init__(self, ):
20          super(Net_C, self).__init__()
21          self.fc1 = nn.Linear(N_STATES, 50)
22          self.fc1.weight.data.normal_(0, 0.1)
23          self.fc2 = nn.Linear(N_ACTIONS, 50)
24          self.fc2.weight.data.normal_(0, 0.1)
25          self.out = nn.Linear(50, 1)
26          self.out.weight.data.normal_(0, 0.1)
27
28      def forward(self, s, a):
29         x = self.fc1(s)
30         y = self.fc2(a)
31         q_value = self.out(F.relu(x+y))
32         return q_value

```

The parameters setting is shown below:

```
1 BATCH_SIZE = 32
2 LR_A = 7.5e-4
3 LR_C = 12e-4
4 GAMMA = 0.9
5 GAMMA_VAR = 0.995
6 TAU = 0.01
7 MEMORY_CAPACITY = 10000
8 env = gym.make('Pendulum-v1')
9 env = env.unwrapped
10 N_ACTIONS = env.action_space.shape[0]
11 N_STATES = env.observation_space.shape[0]
```

4 Analysis of Results

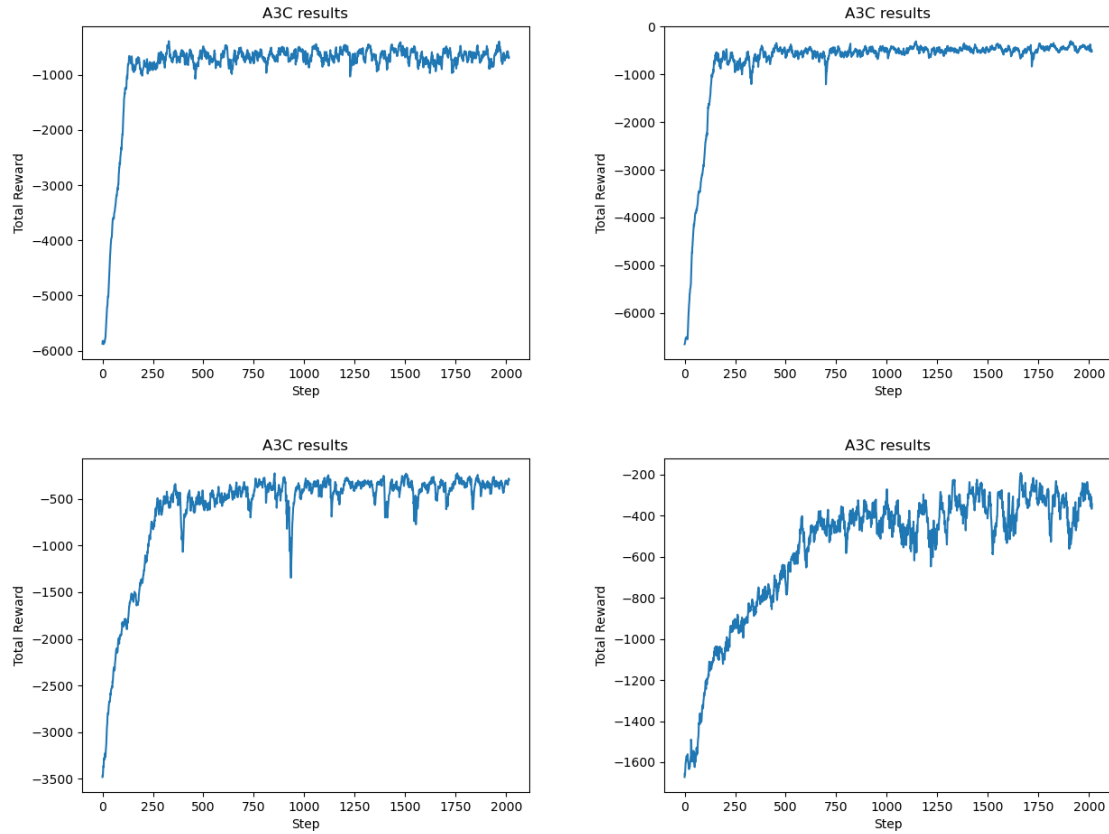
In this section, we first look into the experimental results achieved by A3C algorithm. We will investigate the effect of the maximum number of steps within a single *episode* *MAX_EP_STEP* and the number of *worker* threads on the convergence results of the model using the control variables approach, respectively.

4.1 The result of A3C

4.1.1 Maximum number of steps within an *Episode* *MAX_EP_STEP*

We design 4 groups of experiments to investigate the effect of *MAX_EP_STEP* on the speed and effectiveness of model convergence.

- Group 1: *N_WORKERS* threads 16, *MAX_GLOBAL_EP* is 2000, *MAX_EP_STEP* is 1600
- Group 2: *N_WORKERS* threads being 16, *MAX_GLOBAL_EP* is 2000, *MAX_EP_STEP* is 800
- Group 3: *N_WORKERS* threads 16, *MAX_GLOBAL_EP* is 2000, *MAX_EP_STEP* is 400
- Group 4: *N_WORKERS* threads 16, *MAX_GLOBAL_EP* is 2000, *MAX_EP_STEP* is 200



The upper four results shown from upper left corner to lower right corner is Group1 ~ Group4 respectively.

A comparison of the four results shows that when MAX_EP_STEP is large enough, an appropriate decrease does not cause a large decrease in convergence speed. Such as 200 in this figure is enough to converge to a stable model. However, the larger MAX_EP_STEP is, the more stable the model is when it reaches or converges to convergence. But when the MAX_EP_STEP reaches a certain number, more MAX_EP_STEP will not further improve the convergence and stability, such as 800 in this picture. But when MAX_EP_STEP is lower than this threshold, more MAX_EP_STEP means faster convergence and more stability, and also **much longer training time**. The reason for the above phenomenon may be that the larger MAX_EP_STEP is, the more experience each *worker* has accumulated during the interaction, the more comprehensive and stable the information fed back to the *global network* will be, and accordingly, the less likely to produce large fluctuations. But when MAX_EP_STEP is larger than some threshold, the information learnt becomes redundant and unnecessary.

Summing up, we can get the following conclusions.

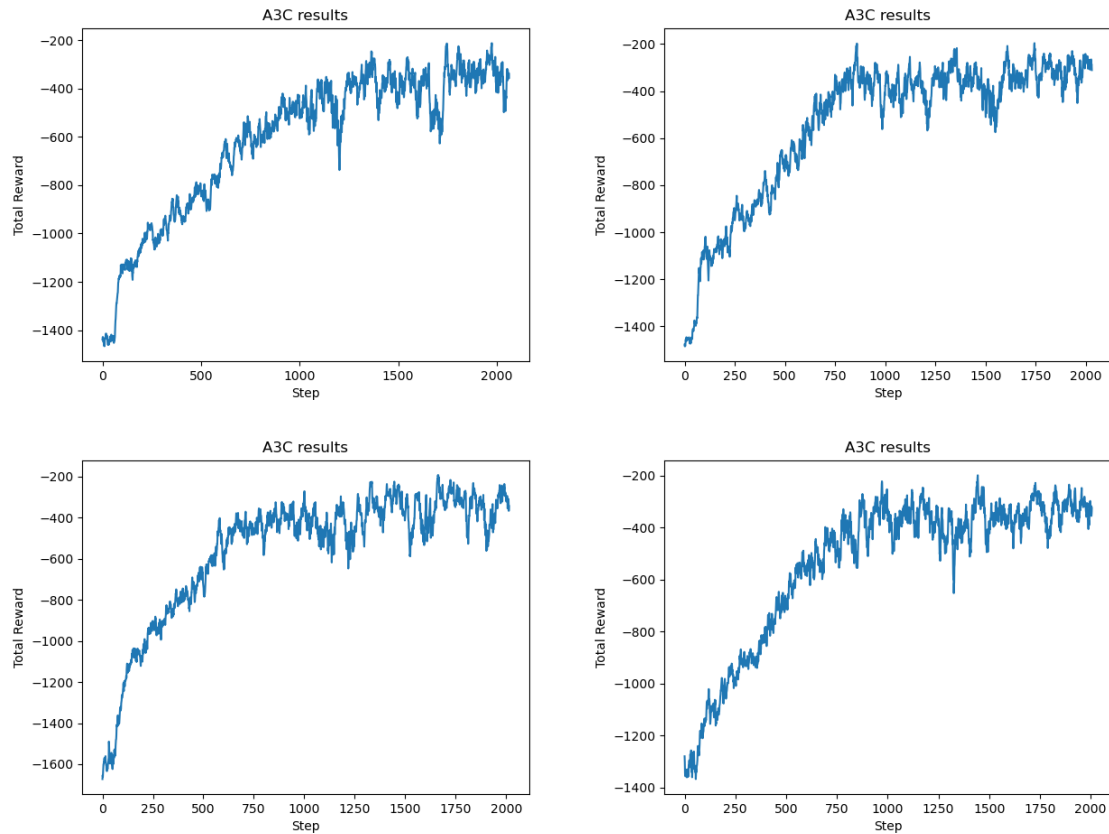
- When MAX_EP_STEP is small, increasing the maximum number of steps within a single *episode* can significantly improve the convergence speed and convergence effect of the model.
- When MAX_EP_STEP is large, increasing the maximum number of steps within a single *episode* will not significantly improve stability and convergence, but leading to redundant training and extra slow training time.

4.1.2 Number of *Worker* threads

We design 4 groups of experiments to investigate the effect of the number of *worker* threads on the convergence speed and effectiveness of the model.

- Group 1: $N_WORKERS$ threads being 64, MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 200
- Group 2: $N_WORKERS$ threads being 32, MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 200
- Group 3: $N_WORKERS$ threads being 16, MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 200

- *Group 4: $N_WORKERS$ threads being 8, MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 200*



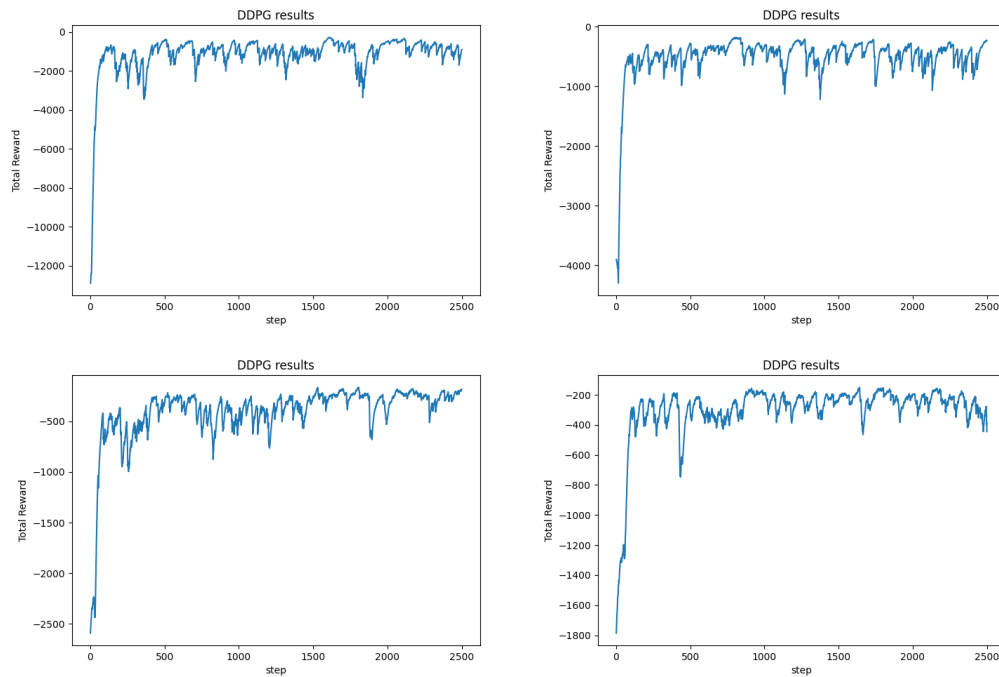
The upper four results shown from upper left corner to lower right corner is Group1 ~ Group4 respectively.

From these results we can see that the number of threads has little effect on training results such as convergence speed and final outcome.

4.2 The result of DDPG

As same as the A3C, we design 4 groups of experiments to investigate the effect of MAX_EP_STEP on the speed and effectiveness of model convergence. Due to the little impact in the A3C experiment, and the limitations of the computer equipment, we do not investigate the effect of the number of *worker* threads.

- *Group 1: MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 1600*
- *Group 2: MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 800*
- *Group 3: MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 400*
- *Group 4: MAX_GLOBAL_EP is 2000, MAX_EP_STEP is 200*



The upper four results shown from upper left corner to lower right corner is Group1 ~ Group4 respectively. First of all, it must be stated that the more *MAX_EP_STEP* is, the longer training time it takes. When *MAX_EP_STEP* is equal to 1600, it takes several times longer than that *MAX_EP_STEP* is equal to 200. In addition we can see that when *MAX_EP_STEP* is equal to 400, stable models can already be trained, and the stability of the model is no worse or even better than the other three models. So *MAX_EP_STEP* should not be set too large to ensure the convergence speed and model stability.

5 Experimental Insights

In this experiment, we first learned about the AC algorithm and the A3C algorithm systematically based on the classroom lectures, and tried out their implementation methods in real practice at the code level. In the process of implementation, we encountered many difficulties due to the difference between pseudo-code to *python* code, including the lack of theoretical basis for debugging the network structure and changing hyperparameters for the best in the process of network design, so I had to rely on manual debugging and compare the results manually. At the same time, since I was not familiar with the use of *tensorflow* tool, I am not quite proficient in class and function designs based on theoretical knowledge. After I consulted a lot of related tutorials I managed to solve this problem.

After that we learned the *DDPG* algorithm based on the classroom lectures. Similar to A3C, in the process of implementing the *DDPG* algorithm, I encountered a number of difficulties such as debugging the network structure and changing hyperparameters. Different with the implement of A3C, the *DDPG* is implemented by *pytorch* which I am not familiar enough. By searching and reading the documentation, I got a better understanding of *DDPG* and *pytorch*.

Solving these problems not only enhances me of my knowledge for A3C and *DDPG*, but also improves my skills in coding, debugging and so on. I learnt more about A3C and *DDPG* methods in reinforcement learning after this project.

