# Homework #2

**Jeongwon Bae (945397461)**

**Abstract**

This project implements the Asynchronous Advantage Actor-Critic (A3C) algorithm using PyTorch to solve the CartPole-v1 environment from OpenAI Gym. The Actor-Critic framework is a fundamental approach in reinforcement learning that combines both policy-based and value-based methods to achieve efficient and stable learning [1]. In this framework, the actor represents the policy $\pi(a_t|s_t; \theta)$, which selects actions $a_t$ based on the current state $s_t$ and policy parameters $\theta$. The critic estimates the value function $V(s_t; \theta_v)$, evaluating the expected return from state $s_t$ using value function parameters $\theta_v$. The actor updates its policy in the direction suggested by the critic's feedback, enabling the agent to learn optimal behaviors.

Building upon this foundation, the Advantage Actor-Critic (A2C) algorithm introduces a synchronous approach where multiple agents interact with their own copies of the environment and update a shared global model in a coordinated manner [2]. In A2C, agents ollect experiences and wait for each other to complete their interactions before synchronously updating the global parameters. This method reduces variance in policy gradient estimates by utilizing the advantage function $A(s_t, a_t)$, defined as:

$$A(s_t, a_t) = Q(s_t, a_t) - V(s_t),$$

where $Q(s_t, a_t)$ is the action-value function, and $V(s_t)$ is the state-value function. The advantage function measures how much better an action $a_t$ is compared to the average expected return from state $s_t$.

The Asynchronous Advantage Actor-Critic (A3C) algorithm extends this approach by allowing agents to operate asynchronously [2]. In A3C, multiple worker processes interact with their own instances of the environment and update the global model independently without waiting for synchronization. This asynchrony offers several benefits: it decorrelates experiences between agents, reducing variance in updates and improving learning stability, and it maximizes computational resource utilization by eliminating idle times due to synchronization.

In the A3C algorithm, each worker maintains a local copy of the global model's parameters. Workers interact with their environments and periodically update the global model based on their own experiences. Each worker collects experience by interacting with its environment over a sequence of time steps. At each time step $t$, the worker observes the current state $s_t$ and selects an action $a_t$ according to its policy $\pi(a_t|s_t; \theta)$. After executing the action, the worker receives a reward $r_t$ and transitions to the next state $s_{t+1}$. This process continues for a fixed number of steps or until a terminal state is reached.

After collecting sufficient experience, the worker computes the n-step return $R_t$ for each time step, which is used to estimate the advantage function $A(s_t, a_t)$. The n-step return is calculated using the rewards and the value estimates from the critic:

$$R_t = \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n V(s_{t+n}; \theta_v),$$

where $\gamma$ is the discount factor (typically $0 < \gamma < 1$), $r_{t+k}$ is the reward received at time $t + k$, and $V(s_{t+n}; \theta_v)$ is the estimated value of the state $s_{t+n}$. If the episode terminates before $n$ steps, the value function term is omitted, and $R_t$ is simply the sum of discounted rewards.

The advantage for each time step is then computed as:

$$A(s_t, a_t) = R_t - V(s_t; \theta_v),$$

where $V(s_t; \theta_v)$ is the critic's estimate of the value of state $s_t$. The advantage function represents how much better or worse the action $a_t$ performed compared to the critic's expected value of the state.

The worker then computes the gradients of the policy and value function losses. The policy loss at each time step is given by:

$$L_{\text{policy}} = -\log \pi(a_t | s_t; \theta) A(s_t, a_t),$$

which encourages the policy to assign higher probabilities to actions that yield positive advantages and lower probabilities to those with negative advantages. This effectively reinforces actions that lead to better-than-expected outcomes.

The value loss is defined as:

$$L_v = \left(R_t - V(s_t; \theta_v)\right)^2,$$

which trains the critic to minimize the difference between the estimated value and the actual return observed. An entropy regularization term is often included to encourage exploration by preventing the policy from becoming too deterministic too quickly:

$$L_{\text{entropy}} = -\beta \sum_a \pi(a | s_t; \theta) \log \pi(a | s_t; \theta),$$

where $\beta$ is a hyperparameter controlling the strength of the entropy regularization.

The total loss for the worker combines these components:

$$L = L_{\text{policy}} + c_v L_v + L_{\text{entropy}},$$

where $c_v$ is a hyperparameter weighting the value loss relative to the policy loss.

The worker computes gradients of the total loss with respect to its local parameters $\theta$ and $\theta_v$. These gradients are then applied to update the global parameters using an optimizer such as RMSProp or Adam. The updates are performed asynchronously, meaning that each worker updates the global model independently without waiting for other workers. After updating the global model, the worker synchronizes its local model parameters with the updated global parameters.

The CartPole-v1 environment is a classic control problem where the objective is to prevent a pole, attached by an unactuated joint to a cart, from falling over by applying forces to move the cart left or right [3]. The environment provides a continuous four-dimensional state space representing the cart position ($x$), cart velocity ($\dot{x}$), pole angle ($\theta$), and pole angular velocity ($\dot{\theta}$). The action

space is discrete, consisting of two actions: pushing the cart to the left (0) or to the right (1). The agent receives a reward of +1 for every time step the pole remains upright. The episode terminates if the pole deviates more than 12° from vertical, the cart moves beyond certain boundaries, or the episode length reaches 500 time steps.

By leveraging the A3C algorithm to train an agent on the CartPole-v1 task, efficient exploration and faster convergence can be achieved compared to synchronous methods. Utilizing multiple asynchronous worker processes, each interacting with its own environment instance, allows for continuous updating of the shared global model based on the experiences of all workers. This approach enables the agent to quickly learn a robust policy that balances the pole for extended periods.

This implementation demonstrates the practical application of the A3C algorithm with clear formulations and parameter updates. It showcases the benefits of asynchronous updates in reducing training time and improving performance. Applying the A3C algorithm to the CartPole-v1 environment validates its effectiveness on a well-known control problem and highlights the advantages of using asynchronous, parallelized training methods in reinforcement learning.

**Step-by-Step Code Explanation**

We will now examine the core components of the reinforcement learning implementation, focusing specifically on the A3C algorithm applied to the CartPole-v1 environment [4]. The code is presented in logical blocks, with detailed explanations and line-by-line comments to demonstrate step-by-step understanding.

Firstly, the necessary libraries are imported, and the hyperparameters are set.

```
1   import gym                                # OpenAI Gym for the environment
2   import torch                              # PyTorch for neural networks
3   import torch.nn as nn                     # neural network modules
4   import torch.nn.functional as F           # functional modules for activation functions
5   import torch.optim as optim               # optimization algorithms
6   from torch.distributions import Categorical # Categorical distribution for action selection
7   import torch.multiprocessing as mp        # multiprocessing for parallel processing
8   import time                               # time module for delays during testing
9
10  # Hyperparameters
11  n_train_processes = 3    # Number of parallel training processes
12  learning_rate = 0.0002   # Learning rate for the optimizer
13  update_interval = 5      # Interval for updating the global model
14  gamma = 0.98             # Discount factor for future rewards
15  max_train_ep = 300       # Maximum number of training episodes per worker
16  max_test_ep = 400        # Maximum number of testing episodes
```

In this block, we import essential modules such as `gym` for the environment, `torch` and its submodules for neural network implementation, and `torch.multiprocessing` for parallel processing. The hyperparameters are defined to control various aspects of the training process, such as the learning rate, the number of episodes, and the discount factor.

Next, the Actor-Critic neural network model is defined.

```
1   class ActorCritic(nn.Module):
2     def __init__(self):
3       super(ActorCritic, self).__init__()
4       self.fc1 = nn.Linear(4, 256)          # Input layer with 4 inputs (state dimensions) and
                256 outputs
5       self.fc_pi = nn.Linear(256, 2)        # Policy output layer with 2 actions (left or right)
6       self.fc_v = nn.Linear(256, 1)         # Value output layer producing a single value
                estimate
7
8     def pi(self, x, softmax_dim=0):
9       x = F.relu(self.fc1(x))               # Apply ReLU activation after first layer
```

```
10          x = self.fc_pi(x)                          # Compute action logits
11          prob = F.softmax(x, dim=softmax_dim)   # Convert logits to probabilities using softmax
12          return prob
13
14      def v(self, x):
15          x = F.relu(self.fc1(x))     # Apply ReLU activation after first layer
16          v = self.fc_v(x)             # Compute state value estimate
17          return v
```

Here, the `ActorCritic` class inherits from `nn.Module` and defines the neural network architecture. The network consists of a shared fully connected layer (`fc1`) followed by separate layers for the policy (`fc_pi`) and value (`fc_v`) outputs. The `pi` method computes the action probabilities using the softmax function, while the `v` method computes the state value estimate.

The `train` function is central to the training process of each worker in the A3C algorithm. Each worker executes this function independently, interacting with its own environment instance and contributing to the asynchronous updates of the global model. It begins by creating a local instance of the `ActorCritic` model, which is synchronized with the global model's parameters. This synchronization ensures that all workers start with the same initial parameters. Additionally, an optimizer is defined to update the global model's parameters, and each worker creates its own instance of the CartPole-v1 environment, allowing them to interact independently.

```
1   def train(global_model, rank):
2       local_model = ActorCritic()        # Create a local model for the worker
3       local_model.load_state_dict(global_model.state_dict()) # Synchronize local model with global
            model
4
5       optimizer = optim.Adam(global_model.parameters(),
6                   lr=learning_rate)  # Define optimizer for global model
7
8       env = gym.make('CartPole-v1')      # Create an instance of the environment
```

Once the environment is set up, the worker enters a loop where it will run for a specified number of episodes. At the start of each episode, the environment is reset to its initial state. The `done` flag is initialized to `False`, indicating that the episode has not yet terminated.

```
1       for n_epi in range(max_train_ep):    # Loop over episodes
2           done = False
3           s = env.reset()                   # Reset the environment to the initial state
```

As the episode progresses, the worker interacts with the environment by collecting experiences. For each time step, the current state is passed through the local policy network to obtain action probabilities. A categorical distribution is created from these probabilities, and an action is sampled. The worker then executes the action in the environment, receiving the next state, reward, and a `done` flag that indicates whether the episode has ended. The current state, action, and scaled reward are stored in lists for later use. This continues until the worker reaches the specified `update_interval` or the episode ends.

```
1       while not done:
2           s_lst, a_lst, r_lst = [], [], []                        # Lists to store states, actions, and
                rewards
3           for t in range(update_interval):
4               prob = local_model.pi(torch.from_numpy(s).float()) # Get action probabilities from policy
                    network
5               m = Categorical(prob)      # Create a categorical distribution based on probabilities
6               a = m.sample().item()      # Sample an action from the distribution
7               s_prime, r, done, info = env.step(a)          # Execute the action in the environment
8
9               s_lst.append(s)             # Store the current state
10              a_lst.append([a])           # Store the action taken
11              r_lst.append(r / 100.0)     # Store the scaled reward (divided by 100)
12
```

```
13        s = s_prime                   # Update the current state
14        if done:
15          break                       # Exit the loop if the episode is done
```

Once the worker has collected a sequence of experiences, it proceeds to compute the necessary values for updating the global model. If the episode has ended, the bootstrap value `R` is set to `0.0`, since there are no future rewards. Otherwise, it is set to the estimated value of the final state, computed by the local critic network. The worker then computes the discounted returns in reverse order from the collected rewards and stores them in a list. This list is reversed to align the returns with the chronological order of the experiences.

```
1        s_final = torch.tensor(s_prime, dtype=torch.float)
2        R = 0.0 if done else local_model.v(s_final).item()   # Compute the bootstrap value
3        td_target_lst = []
4        for reward in r_lst[::-1]:
5          R = reward + gamma * R      # Compute the discounted return
6          td_target_lst.append([R])   # Store the return
7        td_target_lst.reverse()         # Reverse to match the order of experiences
```

Next, the worker prepares the data for updating the model by converting the lists of states, actions, and target returns into PyTorch tensors. The advantage estimates are computed by subtracting the value estimates from the target returns. These advantages represent how much better or worse an action performed compared to the expected outcome.

```
1        # Convert lists to tensors
2        s_batch = torch.tensor(s_lst, dtype=torch.float)
3        a_batch = torch.tensor(a_lst)
4        td_target = torch.tensor(td_target_lst)
5
6        advantage = td_target - local_model.v(s_batch)   # Compute the advantage estimates
```

The loss is then computed, which consists of both a policy loss and a value loss. The policy loss encourages the model to take actions that have higher advantages by minimizing the negative log probabilities of actions weighted by their advantages. The value loss minimizes the difference between the predicted state values and the computed returns using the smooth L1 loss. (We note that the entropy term is missing here, and to fully implement the algorithm as described in the abstract, it should be added.) Gradients are cleared, and the loss is backpropagated through the network to compute the gradients with respect to the model's parameters.

```
1        pi = local_model.pi(s_batch, softmax_dim=1)       # Get action probabilities for the batch
             of states
2        pi_a = pi.gather(1, a_batch)                       # Gather probabilities of the actions taken
3        loss = -torch.log(pi_a) * advantage.detach() + \
4          F.smooth_l1_loss(local_model.v(s_batch), td_target.detach())  # Compute the policy and
             value losses
5        # entropy = - (pi * torch.log(pi + 1e-20)).sum(dim=1)  # Compute entropy (to be added)
6        # loss -= beta * entropy                            # Subtract entropy term from loss
7
8        optimizer.zero_grad()                             # Clear gradients
9        loss.mean().backward()                            # Backpropagate the mean loss
```

Once the gradients are computed, they are manually copied from the local model to the global model's parameters. This step is crucial for ensuring that the global model is updated with the gradients calculated from the local experiences. After copying the gradients, the optimizer is used to update the global model's parameters. Finally, the local model is synchronized again with the updated global model to ensure it has the latest parameters. At the end of the training loop, the environment is closed, and a message is printed to indicate that the worker has completed its training episodes.

```
1        # Copy gradients from local model to global model
```

```
2        for global_param, local_param in zip(global_model.parameters(), local_model.parameters()):
3          global_param._grad = local_param.grad
4        optimizer.step()                              # Update global model parameters
5        local_model.load_state_dict(global_model.state_dict())  # Synchronize local model with
             global model
6    env.close()
7    print(f"Training process {rank} reached maximum episode.")
```

This explanation outlines how the worker collects experiences from the environment, computes the necessary values for updating the model, and asynchronously updates the shared global model. Through this process, the worker contributes to the overall learning of the A3C algorithm, all while interacting with its own instance of the environment independently from other workers.

The `test` function is designed to evaluate the performance of the trained global model. It runs in a separate process and interacts with the environment without updating the model's parameters. This allows us to monitor how well the agent is performing over time.

```
1    def test(global_model):
2      env = gym.make('CartPole-v1')
3      score = 0.0
4      print_interval = 20
5
6      for n_epi in range(max_test_ep):  # Loop over episodes
7        done = False
8        s = env.reset()    # Reset the environment
9        while not done:
10         prob = global_model.pi(torch.from_numpy(s).float())   # Get action probabilities from
                global model
11         a = Categorical(prob).sample().item() # Sample an action
12         s_prime, r, done, info = env.step(a)  # Execute the action
13         s = s_prime    # Update the current state
14         score += r     # Accumulate the reward
15
16       if n_epi % print_interval == 0 and n_epi != 0:
17         print(f"# of episode :{n_epi}, avg score : {score / print_interval:.1f}")
18         score = 0.0
19         time.sleep(1) # Pause for readability
20       env.close()
```

In this function, an instance of the CartPole-v1 environment is created using `gym.make('CartPole-v1')`. The variable `score` is initialized to accumulate rewards, and `print_interval` is set to determine how often the average score is displayed. The function runs a loop over a specified number of episodes (`max_test_ep`). At the start of each episode, the environment is reset using `env.reset()`, and the `done` flag is set to `False`.

Within each episode, the agent interacts with the environment until the episode ends. The current state `s` is converted to a PyTorch tensor and passed through the global policy network `global_model.pi` to obtain action probabilities. A categorical distribution is created from these probabilities, and an action `a` is sampled using `Categorical(prob).sample().item()`. The action is executed in the environment with `env.step(a)`, returning the next state `s_prime`, the reward `r`, and the `done` flag. The current state `s` is updated to `s_prime`, and the reward `r` is added to the cumulative `score`.

After each episode, the function checks if the episode number is a multiple of `print_interval` and not zero. If so, it prints the average score over the last `print_interval` episodes using:

```
1    print(f"# of episode :{n_epi}, avg score : {score / print_interval:.1f}")
```

The cumulative `score` is then reset to `0.0`, and a short pause is added with `time.sleep(1)` to improve readability of the output. Finally, the environment is closed with `env.close()` after all episodes are completed.

Moving on to the main execution block, this is where the global model is initialized, and the training and testing processes are started using multiprocessing.

```
1   if __name__ == '__main__':
2     global_model = ActorCritic()   # Initialize the global model
3     global_model.share_memory()    # Share model parameters for multiprocessing
4
5     processes = []
6     for rank in range(n_train_processes + 1): # Create processes for training and testing
7       if rank == 0:
8         p = mp.Process(target=test, args=(global_model,))    # Create the testing process
9       else:
10        p = mp.Process(target=train, args=(global_model, rank,))  # Create training processes
11      p.start()                    # Start the process
12      processes.append(p)
13    for p in processes:
14      p.join()                     # Wait for all processes to finish
```

In this block, the global model is instantiated by creating an instance of the `ActorCritic` class. The `share_memory()` method is called to allow the model's parameters to be shared across multiple processes, which is essential for the multiprocessing setup. An empty list `processes` is initialized to keep track of all the processes.

A loop is used to create both training and testing processes. If the `rank` is zero, a testing process is created by targeting the `test` function and passing the `global_model` as an argument. For other ranks, training processes are created by targeting the `train` function and passing the `global_model` and the worker's `rank`. Each process is started with `p.start()`, and the process object is appended to the `processes` list.

After starting all processes, another loop is used to call `p.join()` on each process, causing the main program to wait for all child processes to complete before exiting. This ensures that all training and testing operations are finished properly.

By organizing the code in this way, the program leverages multiprocessing to train multiple agents in parallel while also evaluating the agent's performance simultaneously. This setup is fundamental to the A3C algorithm, where asynchronous updates from multiple workers contribute to learning a shared policy and value function.

Overall, the code implements an A3C agent that learns to solve the CartPole-v1 environment. The `ActorCritic` model defines the neural network architecture, combining both policy and value outputs. The `train` function handles experience collection, computation of advantages, and asynchronous updates to the global model. The `test` function evaluates the agent's performance by running it in the environment without learning. The main execution block sets up the multiprocessing environment, initializes the global model, and starts the training and testing processes.

This implementation demonstrates how asynchronous methods can effectively utilize computational resources and lead to efficient learning in reinforcement learning tasks. By having multiple worker processes interact with their own environments and update a shared model, the agent can learn robust policies that generalize well across different states and scenarios.

### Complete Code Script

```python
1    import gym                                      # Import OpenAI Gym for the environment
2    import torch                                    # Import PyTorch library
3    import torch.nn as nn                           # Import neural network modules
4    import torch.nn.functional as F                 # Import functional modules for activation functions
5    import torch.optim as optim                     # Import optimization algorithms
6    from torch.distributions import Categorical # Import Categorical distribution for probabilistic
         action selection
7    import torch.multiprocessing as mp             # Import multiprocessing for parallel execution
8    import time                                     # Import time module for delays during testing
9
10   # Hyperparameters
11   n_train_processes = 3    # Number of parallel training processes
12   learning_rate = 0.0002   # Learning rate for optimizer
13   update_interval = 5      # Interval for updating global model
14   gamma = 0.98             # Discount factor for future rewards
15   max_train_ep = 300       # Maximum number of training episodes per worker
16   max_test_ep = 400        # Maximum number of testing episodes
17
18   class ActorCritic(nn.Module):
19     def __init__(self):
20       super(ActorCritic, self).__init__()
21       self.fc1 = nn.Linear(4, 256)           # Input layer to hidden layer with 256 units
22       self.fc_pi = nn.Linear(256, 2)         # Output layer for policy (2 possible actions)
23       self.fc_v = nn.Linear(256, 1)          # Output layer for value function
24
25     def pi(self, x, softmax_dim=0):
26       x = F.relu(self.fc1(x))                # Apply ReLU activation function
27       x = self.fc_pi(x)                      # Pass through policy output layer
28       prob = F.softmax(x, dim=softmax_dim)   # Convert logits to probabilities with softmax
29       return prob
30
31     def v(self, x):
32       x = F.relu(self.fc1(x))                # Apply ReLU activation function
33       v = self.fc_v(x)                       # Pass through value output layer
34       return v
35
36   def train(global_model, rank):
37     local_model = ActorCritic()                           # Create a local model
38     local_model.load_state_dict(global_model.state_dict())   # Synchronize with global model
39
40     optimizer = optim.Adam(global_model.parameters(), lr=learning_rate)  # Optimizer for global
         model
41
42     env = gym.make('CartPole-v1')                          # Initialize environment
43
44     for n_epi in range(max_train_ep):                      # Training loop over episodes
45       done = False
46       s = env.reset()                                      # Reset environment
47       while not done:
48         s_lst, a_lst, r_lst = [], [], []                   # Lists to store states, actions, rewards
49         for t in range(update_interval):
50           prob = local_model.pi(torch.from_numpy(s).float())  # Get action probabilities
51           m = Categorical(prob)                            # Create categorical distribution
52           a = m.sample().item()                            # Sample action
53           s_prime, r, done, info = env.step(a)             # Perform action in environment
54
55           s_lst.append(s)                                  # Append state to list
56           a_lst.append([a])                                # Append action to list
57           r_lst.append(r / 100.0)                          # Append scaled reward to list
58
59           s = s_prime                                      # Update state
60           if done:
61             break                                          # Break if episode is done
62
63         s_final = torch.tensor(s_prime, dtype=torch.float)
64         R = 0.0 if done else local_model.v(s_final).item()   # Bootstrap value for last state
65         td_target_lst = []
66         for reward in r_lst[::-1]:
67           R = reward + gamma * R                           # Calculate discounted return
68           td_target_lst.append([R])                        # Append to target list
```

```
69         td_target_lst.reverse()                            # Reverse to match order
70
71         s_batch = torch.tensor(s_lst, dtype=torch.float) # Convert states to tensor
72         a_batch = torch.tensor(a_lst)                      # Convert actions to tensor
73         td_target = torch.tensor(td_target_lst)            # Convert targets to tensor
74
75         advantage = td_target - local_model.v(s_batch)     # Compute advantage
76
77         pi = local_model.pi(s_batch, softmax_dim=1)        # Get action probabilities
78         pi_a = pi.gather(1, a_batch)                       # Get probabilities of taken actions
79         loss = -torch.log(pi_a) * advantage.detach() + \
80           F.smooth_l1_loss(local_model.v(s_batch), td_target.detach())  # Compute loss
81         # entropy = - (pi * torch.log(pi + 1e-20)).sum(dim=1)  # Compute entropy (to be added)
82         # loss -= beta * entropy                            # Subtract entropy term from loss
83
84         optimizer.zero_grad()                              # Clear gradients
85         loss.mean().backward()                             # Backpropagate loss
86         for global_param, local_param in zip(global_model.parameters(), local_model.parameters()):
87           global_param._grad = local_param.grad         # Copy gradients to global model
88         optimizer.step()                                   # Update global model
89         local_model.load_state_dict(global_model.state_dict())  # Synchronize local model
90
91     env.close()
92     print(f"Training process {rank} reached maximum episode.")
93
94   def test(global_model):
95     env = gym.make('CartPole-v1')                          # Initialize environment
96     score = 0.0
97     print_interval = 20
98
99     for n_epi in range(max_test_ep):                       # Testing loop over episodes
100       done = False
101       s = env.reset()                                      # Reset environment
102       while not done:
103         prob = global_model.pi(torch.from_numpy(s).float())   # Get action probabilities
104         a = Categorical(prob).sample().item()            # Sample action
105         s_prime, r, done, info = env.step(a)             # Perform action
106         s = s_prime                                       # Update state
107         score += r                                        # Accumulate reward
108
109       if n_epi % print_interval == 0 and n_epi != 0:
110         print(f"# of episode :{n_epi}, avg score : {score / print_interval:.1f}")
111         score = 0.0
112         time.sleep(1)                                     # Pause for readability
113     env.close()
114
115   if __name__ == '__main__':
116     global_model = ActorCritic()                       # Initialize global model
117     global_model.share_memory()                          # Share model parameters for multiprocessing
118
119     processes = []
120     for rank in range(n_train_processes + 1):          # Create processes
121       if rank == 0:
122         p = mp.Process(target=test, args=(global_model,))    # Testing process
123       else:
124         p = mp.Process(target=train, args=(global_model, rank,))  # Training process
125       p.start()                                          # Start process
126       processes.append(p)
127     for p in processes:
128       p.join()                                           # Wait for all processes to finish
```

# References

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, Second. The MIT Press, 2018.

[2] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proceedings of The 33rd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 48, New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1928–1937.

[3] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuronlike adaptive elements that can solve difficult learning control problems," in *Artificial Neural Networks: Concept Learning*. IEEE Press, 1990, pp. 81–93, ISBN: 0818620153.

[4] S. Rho. "Minimalrl-pytorch. "[Online]. Available: `https://github.com/seungeunrho/minimalRL`.