

Bonus Point Assignment II

Submission Evaluation Group 40

Overview

Question	Points
A2	14/14
B1	4/4
B2	13/14
C1	12/12
Total	43/44
Bonus points	6/6

6 bonus points in this assignment correspond to 5% of bonus for the final exam.

Below follows the detailed evaluation of your submission. Red boxes are errors which deducted points. In case of unittest failures, code is given with which you can reproduce the error. Yellow boxes denote comments for your information.

A2

a) Defining the Policy

```
1 class Policy(nn.Module):
2     def __init__(self):
3         super(Policy, self).__init__()
4         self.fully_connected_layer = nn.Sequential(
5             nn.Linear(4, 128), nn.ReLU(), nn.Linear(128, 128), nn.ReLU(),
6             nn.Linear(128, 2)
7         )
8         self.softmax_layer = nn.Softmax(1)
9
10    def forward(self, x):
11        x = self.fully_connected_layer(x)
12        x = self.softmax_layer(x)
13        return x
```

ok

Network uses 17,410 weights.

b) Action Sampling

```
1 def sample_action(probs):
2     action_distribution = torch.distributions.Categorical(probs=probs)
3     action = action_distribution.sample()
4     log_prob = action_distribution.log_prob(action)
5     return action, log_prob
```

ok

c) Estimate Return

```
1 def estimate_return(rewards, gamma=0.99):
2     returns = []
3     reward = 0
4     for r in reversed(rewards):
5         reward = reward * gamma + r
6         returns.insert(0, reward)
7     mean = np.mean(returns)
8     std = np.std(returns)
9     returns = (returns - mean) / std
10    return returns
```

ok

d) Training Loop

ok

Average validation reward: 500.0

Average test reward: 500.0

B1

a) Implement Replay Buffer

```
15 def add(self, state, action, reward, next_state, is_terminal):
16     self.state_memory[self.ptr] = state
17     self.next_state_memory[self.ptr] = next_state
18     self.action_memory[self.ptr] = action
19     self.reward_memory[self.ptr] = reward
20     self.terminal_memory[self.ptr] = is_terminal
21     self.mem_cntr = self.mem_cntr + 1 if not self.is_filled() else self.mem_size
22     self.ptr = (self.ptr + 1) % self.mem_size
```

ok

```

24 def sample_batch(self, batch_size):
25     index = np.random.choice(range(self.mem_cntr), size=batch_size, replace=False)
26     states = self.state_memory[index]
27     actions = self.action_memory[index]
28     rewards = self.reward_memory[index]
29     next_states = self.next_state_memory[index]
30     is_terminal = self.terminal_memory[index]
31     return states, actions, rewards, next_states, is_terminal

```

ok

b) Fill replay buffer

```

3 state = env.reset()
4 while not buffer.is_filled():
5     action = env.action_space.sample()
6     next_state, reward, done, _ = env.step(action)
7     buffer.add(state, action, reward, next_state, done)
8     if done:
9         state = env.reset()
10    else:
11        state = next_state

```

ok

B2

a) Define Q-Network

```

6 class DeepQNetwork(nn.Module):
7     def __init__(self):
8         super(DeepQNetwork, self).__init__()
9         self.fully_connected_layer = nn.Sequential(
10             nn.Linear(4, 128), nn.ReLU(), nn.Linear(128, 128), nn.ReLU(),
11             nn.Linear(128, 2)
12         )
13
14     def forward(self, state):
15         Q = self.fully_connected_layer(state)
16         return Q

```

ok

Network uses 17,410 weights.

b) ϵ -Greedy

```

1 def epsilon_greedy(state, q_network, epsilon=0.05):
2     state = torch.tensor(state, dtype=torch.float32)
3     if torch.rand(1) < epsilon:
4         action = env.action_space.sample()
5     else:
6         action = q_network(state).argmax().item()
7     return action

```

AssertionError

```

class DummyModule(nn.Module):
    def forward(self, state):
        return torch.FloatTensor([-1, -2])

actions = [
    epsilon_greedy(np.array([1, 2, 3, 4]), DummyModule(), epsilon=0.5)
    for _ in range(1000)
]
action_0 = 1000 - sum(actions)
assert 704 <= action_0 <= 794, f"Frequency {action_0} is outside the 99.9% confidence interval"

```

Frequency 795 is outside the 99.9% confidence interval

c) Loss Function

```

4 def compute_loss(
5     q_network,
6     target_network,
7     states,
8     actions,
9     rewards,
10    next_states,
11    is_terminal,
12    gamma=0.99
13 ):
14     qvals = torch.gather(q_network(states), 1, actions)
15     max_actions = q_network(next_states).argmax(dim=1, keepdim=True)
16     expected_qvals = torch.gather(target_network(next_states), 1, max_actions)
17     expected_qvals = expected_qvals * gamma * torch.logical_not(is_terminal) + rewards
18     loss = mse(qvals, expected_qvals.detach())
19     return loss

```

ok

d) Training

ok

Average validation reward: 500.0

Average test reward: 500.0

C1

a) Computing MC estimates

```
18 for reward, is_terminal in zip(
19     reversed(memory.rewards), reversed(memory.is_terminals)
20 ):
21     if not is_terminal:
22         discounted_reward = discounted_reward * self.gamma + reward
23     else:
24         discounted_reward = reward
```

ok

b) Computing surrogate loss functions

```
36 advantages = rewards - state_values.detach()
37 surr1 = ratios * advantages
38 surr2 = torch.clamp(
39     ratios, min=1 - self.eps_clip, max=1 + self.eps_clip
40 ) * advantages
```

ok

c) Data Collection

```
30 for t in range(max_timesteps):
31     timestep += 1
32     action = ppo.policy_old.act(state, memory)
33     next_state, reward, done, _ = env.step(action)
34     memory.rewards.append(reward)
35     memory.is_terminals.append(done)
36     if done:
37         state = env.reset()
38     else:
39         state = next_state
```

ok

d) Policy Update

```
40 if timestep % update_timestep == 0:  
41     ppo.update(memory)  
42     memory.clear_memory()
```

ok