

Understanding Reinforcement Learning Code

Xin Gu

The reinforcement learning code is chosen from PyTorch Reinforcement Learning (DQN) Tutorial¹.

1 High-level Overview

This code implements a Deep Q-Network (DQN) to train an agent on the CartPole-v1 environment using PyTorch. The DQN algorithm enables the agent to learn an optimal policy by 1) predicting expected rewards for each action and adjusting its predictions through experience, 2) a replay memory stores past experiences, allowing the agent to sample them for training, which stabilizes the learning process, 3) the Q-network predicts action values, while a target network provides stable estimates during updates. Through repeated interaction with the environment, the agent learns to keep the pole balanced by improving its policy iteratively, leading to longer balanced episodes over time.

2 Core Section Comments

```
1 # Function to optimize the Q-network based on a batch of past experiences
2 def optimize_model():
3     # Ensure there's enough data in replay memory for a full batch
4     if len(memory) < BATCH_SIZE:
5         return # Exit if not enough transitions are stored
6
7     # Sample a batch of transitions from replay memory
8     transitions = memory.sample(BATCH_SIZE)
9     # Reshape transitions
10    batch = Transition(*zip(*transitions))
11
12    # Compute a mask of non-final states (next_state is not None)
13    # Allows the model to ignore terminal states in the Q-value calculation
14    non_final_mask = torch.tensor(tuple(map(lambda s: s is not None,
15                                           batch.next_state)), device=device,
16                                dtype=torch.bool)
17    # Concatenate all non-final next states into a single tensor
18    non_final_next_states = torch.cat([s for s in batch.next_state if s is not
19                                     None])
20
21    # Concatenate each component of the batch into tensors for batch processing
22    state_batch = torch.cat(batch.state)
23    action_batch = torch.cat(batch.action)
24    reward_batch = torch.cat(batch.reward)
```

¹https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html#reinforcement-learning-dqn-tutorial

```

23
24     # Compute Q(s_t, a) - the Q-values for the actions taken in each state in the
    batch
25     state_action_values = policy_net(state_batch).gather(1, action_batch)
26
27     # Compute V(s_{t+1}) for all next states using the target network
28     next_state_values = torch.zeros(BATCH_SIZE, device=device)
29     with torch.no_grad():
30         # For non-final states, choose the action with the max Q-value in the
    target network
31         next_state_values[non_final_mask] = target_net(non_final_next_states).max
    (1).values
32
33     # Compute the expected Q values (target values) using the Bellman equation
34     expected_state_action_values = (next_state_values * GAMMA) + reward_batch
35
36     # Compute the loss between the Q-values from policy_net and the target Q-
    values
37     criterion = nn.SmoothL1Loss()
38     loss = criterion(state_action_values, expected_state_action_values.unsqueeze
    (1))
39
40     # Perform a backward pass to compute gradients and update weights
41     optimizer.zero_grad() # Clear any previous gradients
42     loss.backward() # Backpropagate the loss to compute gradients
43
44     # Clip gradients to a max value (100) to prevent exploding gradients
45     torch.nn.utils.clip_grad_value_(policy_net.parameters(), 100)
46     optimizer.step() # Update policy_net's parameters based on gradients

```

Listing 1: One step of optimization

```

1 # Main loop for training across multiple episodes
2 for i_episode in range(num_episodes):
3     # Reset the environment at the beginning of each episode, getting initial
    state
4     state, info = env.reset()
5     # Convert state to a tensor for compatibility with PyTorch
6     state = torch.tensor(state, dtype=torch.float32, device=device).unsqueeze(0)
7
8     for t in count():
9         # Select an action based on the current policy and epsilon-greedy
    exploration
10        action = select_action(state)
11
12        # Take the selected action and observe the next state, reward, and done
    flag
13        observation, reward, terminated, truncated, _ = env.step(action.item())
14        # Convert the reward to a tensor for consistency in optimization step
15        reward = torch.tensor([reward], device=device)
16
17        # Check if the episode is done (pole fell or cart went out of bounds)
18        done = terminated or truncated
19
20        # Set the next state; if the episode is done, next_state is None
21        if terminated:
22            next_state = None
23        else:
24            # Convert next state to tensor

```

```

25         next_state = torch.tensor(observation, dtype=torch.float32, device=
device).unsqueeze(0)
26
27         # Store the current transition (state, action, next_state, reward) in
replay memory
28         memory.push(state, action, next_state, reward)
29
30         # Move to the next state for the next time step
31         state = next_state
32
33         # Perform one optimization step to improve the policy network based on
sampled transitions
34         optimize_model()
35
36         # Soft update of the target network's weights for stability in training
37         # target_net weights = tau * policy_net weights + (1 - tau) * target_net
weights
38         target_net_state_dict = target_net.state_dict()
39         policy_net_state_dict = policy_net.state_dict()
40         for key in policy_net_state_dict:
41             target_net_state_dict[key] = policy_net_state_dict[key]*TAU +
target_net_state_dict[key]*(1-TAU)
42         target_net.load_state_dict(target_net_state_dict)
43
44         # If episode is done, record its duration and break the loop
45         if done:
46             episode_durations.append(t + 1)
47             plot_durations() # Update the plot to show episode durations
48             break

```

Listing 2: Training loop

3 Conclusion

In conclusion, the Rodinia "openmp/hotspot" benchmark showed significant sensitivity to CPU frequency changes, with a significant reduction in execution time as CPU frequency increased. Memory technology also impacted performance, but to a less degree. By analyzing the code, we know that Rodinia "openmp/hotspot" benchmark has a structured data access pattern, meaning good spatial and temporal locality, thus allowing the application to efficiently use the cache, resulting in lower cache miss rates.