Homework 2

Name – Darshan Chudiwal
PSU ID - 926435336

Abstract –

This code implements a simple Deep Q-Network (DQN) algorithm to train an agent for
reinforcement learning tasks. It includes a QNetwork class that defines a neural network to
approximate the action-value function, a ReplayBuffer class for storing and sampling past
experiences to stabilize learning, a DQNAgent class that handles action selection, learns from
experiences, and updates network parameters, and a training loop where the agent interacts
with the environment, collects experiences, and learns to maximize cumulative rewards.
Overall, it demonstrates how to train a DQN agent to learn optimal policies through
interaction with an environment. The environment here could be any environment from the
Gymnasium library, but I wrote the code with LunarLander in mind.

Code – It is something I worked on, inspired by this article –
https://www.katnoria.com/nb_dqn_lunar/

QNetwork class :

```
1. # Define the QNetwork class inheriting from PyTorch's neural network module
2. class QNetwork(nn.Module):
3.     # Initialize the network with state and action sizes and a random seed
4.     def __init__(self, state_size, action_size, seed):
5.         super(QNetwork, self).__init__()
6.         # Set the random seed for reproducibility
7.         self.seed = torch.manual_seed(seed)
8.         # Define the first fully connected layer (input layer). This is a simple
implementation and does not use CNN as described in the DQN paper.
9.         self.fc1 = nn.Linear(state_size, 128)
10.        # Define the second fully connected layer (hidden layer)
11.        self.fc2 = nn.Linear(128, 128)
12.        # Define the final fully connected layer (output layer)
13.        self.fc3 = nn.Linear(128, action_size)
14.
15.    # Define the forward pass through the network
16.    def forward(self, x):
17.        # Pass input through the first layer and apply ReLU activation
18.        x = F.relu(self.fc1(x))
19.        # Pass through the second layer with ReLU activation
20.        x = F.relu(self.fc2(x))
21.        # Pass through the final layer to get action values (Q-values)
22.        return self.fc3(x)
```

Replay Buffer :

```
1. # Define the ReplayBuffer class to store and sample experiences
2. class ReplayBuffer:
3.     # Initialize the buffer with size, batch size, and seed
4.     def __init__(self, buffer_size, batch_size, seed):
5.         self.batch_size = batch_size
6.         # Set the random seed for reproducibility
7.         self.seed = random.seed(seed)
8.         # Initialize the memory as a deque with a maximum length
9.         self.memory = deque(maxlen=buffer_size)
10.        # Define a named tuple to store experiences and later use from
11.        self.experience = namedtuple("Experience", field_names=["observation", "action",
"reward", "next_state", "terminated"])
```

```
12.
13.     # Method to add a new experience to memory
14.     def add(self, observation, action, reward, next_state, terminated):
15.         # Handle cases where observation is a tuple (from certain environments, in this case
the environment was LunarLandingv2)
16.         if isinstance(observation, tuple):
17.             observation = observation[0]
18.         if isinstance(next_state, tuple):
19.             next_state = next_state[0]
20.         # Create an experience tuple
21.         experience = self.experience(observation, action, reward, next_state, terminated)
22.         # Append the experience to the memory buffer
23.         self.memory.append(experience)
24.
25.     # Method to sample a batch of experiences from memory
26.     def sample(self):
27.         # Randomly sample experiences equal to the batch size
28.         experiences = random.sample(self.memory, k=self.batch_size)
29.         # Convert observations to a tensor
30.         observations = torch.from_numpy(np.vstack([experience.observation for experience in
experiences])).float().to(device)
31.         # Convert actions to a tensor and reshape
32.         actions = torch.from_numpy(np.vstack([experience.action for experience in
experiences]).reshape(-1, 1)).long().to(device)
33.         # Convert rewards to a tensor and reshape
34.         rewards = torch.from_numpy(np.vstack([experience.reward for experience in
experiences]).reshape(-1, 1)).float().to(device)
35.         # Convert next states to a tensor
36.         next_states = torch.from_numpy(np.vstack([experience.next_state for experience in
experiences])).float().to(device)
37.         # Convert termination flags to a tensor, cast to uint8, and reshape
38.         terminateds = torch.from_numpy(np.vstack([experience.terminated for experience in
experiences]).astype(np.uint8).reshape(-1, 1)).float().to(device)
39.         # Return the sampled experiences as tensors
40.         return (observations, actions, rewards, next_states, terminateds)
41.
42.     # Return the current size of the memory buffer
43.     def __len__(self):
44.         return len(self.memory)
```

DQNAgent :

```
 1. # Hyperparameters
 2. BUFFER_SIZE = int(1e5)   # Size of the replay memory
 3. BATCH_SIZE = 64          # Number of experiences to sample
 4. GAMMA = 0.99             # Discount factor for future rewards
 5. TAU = 1e-3               # Soft update parameter for the target network
 6. LR = 1e-4                # Learning rate for the optimizer
 7. UPDATE_EVERY = 4         # Frequency of network updates
 8.
 9. # Define the DQNAgent class
10. class DQNAgent:
11.     # Initialize the agent with state and action sizes and a seed
12.     def __init__(self, state_size, action_size, seed):
13.         self.state_size = state_size
14.         self.action_size = action_size
15.         # Set the random seed for reproducibility
16.         self.seed = random.seed(seed)
17.         # Initialize the Q-network (local network)
18.         self.q_network = QNetwork(state_size, action_size, seed).to(device)
19.         # Initialize the fixed Q-network (target network)
20.         self.fixed_network = QNetwork(state_size, action_size, seed).to(device)
21.         # Define the optimizer for the Q-network
22.         self.optimizer = optim.Adam(self.q_network.parameters())
23.         # Initialize the replay memory
24.         self.memory = ReplayBuffer(BUFFER_SIZE, BATCH_SIZE, seed)
25.         # Initialize the timestep counter
26.         self.timestep = 0
27.
```

```python
28.      # Method to process a step in the environment
29.      def step(self, observation, action, reward, next_state, terminated):
30.          # Add the experience to replay memory
31.          self.memory.add(observation, action, reward, next_state, terminated)
32.          # Increment the timestep
33.          self.timestep += 1
34.          # If it's time to update the network
35.          if self.timestep % UPDATE_EVERY == 0:
36.              # Ensure there are enough experiences in memory
37.              if len(self.memory) > BATCH_SIZE:
38.                  # Sample a batch of experiences
39.                  sampled_experiences = self.memory.sample()
40.                  # Learn from the sampled experiences
41.                  self.learn(sampled_experiences)
42.
43.      # Method to learn from a batch of experiences
44.      def learn(self, experiences):
45.          # Unpack the experiences
46.          states, actions, rewards, next_states, terminateds = experiences
47.          # Get the Q-values from the target network for next states (detach to avoid
gradients)
48.          action_values = self.fixed_network(next_states).detach()
49.          # Get the maximum Q-value for each next state
50.          max_action_values = action_values.max(1)[0].unsqueeze(1)
51.          # Compute the target Q-values using the Bellman equation
52.          Q_target = rewards + (GAMMA * max_action_values * (1 - terminateds))
53.          # Get the expected Q-values from the local network
54.          Q_expected = self.q_network(states).gather(1, actions)
55.          # Calculate the loss between expected and target Q-values
56.          loss = F.mse_loss(Q_expected, Q_target)
57.          # Zero the parameter gradients
58.          self.optimizer.zero_grad()
59.          # Perform backpropagation
60.          loss.backward()
61.          # Update the network weights
62.          self.optimizer.step()
63.          # Soft update the target network
64.          self.update_fixed_network(self.q_network, self.fixed_network)
65.
66.      # Method to update the target network parameters
67.      def update_fixed_network(self, q_network, fixed_network):
68.          # Iterate over parameters of both networks
69.          for source_parameters, target_parameters in zip(q_network.parameters(),
fixed_network.parameters()):
70.              # Perform soft update of target network parameters
71.              target_parameters.data.copy_(TAU * source_parameters.data + (1.0 - TAU) *
target_parameters.data)
72.
73.      # Method to select an action using an epsilon-greedy policy
74.      def act(self, observation, eps=0.0):
75.          # Handle cases where observation is a tuple
76.          if isinstance(observation, tuple):
77.              observation = observation[0]
78.          # Generate a random number for epsilon-greedy action selection
79.          rnd = random.random()
80.          # If random number is less than epsilon, select a random action (exploration)
81.          if rnd < eps:
82.              return np.random.randint(self.action_size)
83.          else:
84.              # Convert observation to a tensor and add a batch dimension
85.              observation = torch.from_numpy(observation).float().unsqueeze(0).to(device)
86.              # Set the network to evaluation mode
87.              self.q_network.eval()
88.              # Disable gradient calculation
89.              with torch.no_grad():
90.                  # Get action values from the Q-network
91.                  action_values = self.q_network(observation)
92.              # Set the network back to training mode
93.              self.q_network.train()
94.              # Select the action with the highest Q-value
```

```
 95.                action = np.argmax(action_values.cpu().data.numpy())
 96.                return action
 97.
 98.        # Method to save the trained model weights
 99.        def checkpoint(self, filename):
100.            torch.save(self.q_network.state_dict(), filename)
```

Training Loop :

```
 1. # Instantiate the DQN agent
 2. dqn_agent = DQNAgent(state_size, action_size, seed=0)
 3.
 4. # Initialize a list to keep track of scores
 5. dqn_scores = []
 6. # Initialize a deque to store scores of the last 100 episodes
 7. dqn_scores_window = deque(maxlen=100)
 8. # Initialize epsilon for the epsilon-greedy policy
 9. eps = EPS_START
10. # Record the start time of training
11. start = time()
12. # Initialize a list to store actions taken in each episode
13. dqn_actions = []
14.
15. # Loop over episodes
16. for episode in range(1, MAX_EPISODES + 1):
17.     # Reset the environment and get the initial state
18.     state = env.reset(seed = 42)
19.     # Initialize the score for the current episode
20.     score = 0
21.     # Initialize a list to store actions in the current episode
22.     episode_actions = []
23.
24.     # Loop over steps within an episode
25.     for t in range(MAX_STEPS):
26.         # Select an action using the agent's policy
27.         action = dqn_agent.act(state, eps)
28.         # Append the action to the episode's action list
29.         episode_actions.append(action)
30.         # Take the action in the environment and observe the outcome
31.         observation, reward, terminated, truncated, info = env.step(action)
32.         # Handle cases where observation is a tuple
33.         if isinstance(observation, tuple):
34.             observation = observation[0]
35.         # Let the agent process the step (store experience and learn)
36.         dqn_agent.step(state, action, reward, observation, terminated)
37.         # Update the state to the next state
38.         state = observation
39.         # Accumulate the reward to the score
40.         score += reward
41.         # If the episode is terminated, exit the loop
42.         if terminated:
43.             break
44.         # Decay epsilon
45.         eps = max(eps * EPS_DECAY, EPS_MIN)
46.         mean_score = 0
47.         # Every PRINT_EVERY episodes, print the progress
48.         if episode % PRINT_EVERY == 0:
49.             # Calculate the average score over the last 100 episodes
50.             mean_score = np.mean(dqn_scores_window)
51.             print('\r Progress {}/{}, average score:{:.2f}'.format(episode, MAX_EPISODES,
mean_score), end="")
52.             # If the environment is considered solved
53.             if mean_score >= ENV_SOLVED:
54.                 print('\rEnvironment solved in {} episodes, average score:
{:.2f}'.format(episode, mean_score), end="")
55.                 sys.stdout.flush()
56.                 # Save the trained model weights
57.                 dqn_agent.checkpoint('solved_200.pth')
```

```python
58.              break
59.          # Append the score of the episode to the deque and list
60.          dqn_scores_window.append(score)
61.          dqn_scores.append(score)
62.          # Append the actions taken in the episode to the list
63.          dqn_actions.append(episode_actions)
64.
65.  # Record the end time of training
66.  end = time()
67.  # Print the total training time
68.  print('Took {} seconds'.format(end - start))
69.  # Store the total training time
70.  dqn_time = end-start
```