# Moon Landing using Deep Q-Learning

Khang Pham (A20464881)

March 22, 2025

Course: MMAE 500

Professor: Scott Dawson

# Contents

**Abstract**

This paper discussed the implementation of Deep Q-Learning, a model-free reinforcement learning algorithm on solving continuous state space environments. In contrast to other Data-Driven Modeling method, reinforcement learning create its own data set to learn from with out supervision from experts. To aid application of Deep Q-Learning, other techniques such as epsilon-greedy and soft-updated target network are also introduced. The details on these methods are discussed in section 2. The environment is highly controlled for demonstration purposes; however, the techniques is expected to able to solve higher dimension problems. As a result of the project, the agent successfully learned the environment and found the optimal solution to the problem i.e. land the spaceship fast and accurately. Some other potential strategies to apply Deep Q-Learning are also discussed in sections 2 and 5. The code for this project can be found in the appendix **??**

# 1 Introduction

Traditional control theory techniques aim to describe the systems based on determined rules and correlations between variables. These techniques, with correct calibration, often present accurate results with robust computational potentials. However, with complex systems, sometimes it is difficult to derive governed equations and control them. This project aims to experiment with using reinforcement learning to solve the high-dimension problem: landing a space shuttle.

Q-learning is the original studied method. However, understanding the limitation of discrete state and action space, Deep Q-Learning is implemented. Although reinforcement learning is an unstable method with accuracy often lower than that of control theory techniques, the benefit of model-free control is still appealing. Without the need to study and model the systems in detail, it is possible to solve much more complicated problems, including uncertainty problems.

# 2 Methods

The result discussed in section 4 is achieved by Deep Q-Learning. This method combines Q-Learning (sec. 2.1) and neural networks (sec. 2.2. Details on why the neural network was introduced to improve the method are also discussed in the following section.

## 2.1 Q-Learning

Q-learning is a model-free reinforcement learning algorithm to evaluate the value at a particular state [1]. The algorithm assumes the Markov property, meaning the correlation between the current action and future state depends exclusively on the present state and not on any of the previous states [2].

Q-Learning algorithm can be described in three steps:

1. Initialize the action-state table.
2. Choose an action based on the current state.
3. Update Q-value.

### 2.1.1 Initialize action-state table

The action-state table lists out all action-state pairs with a Q-value assigned to each pair. This step presents a limitation of the Q-Learning algorithm which is the dimensions of possible states and actions are both required to be finite. It is suitable for cases where the state space and action space are both discrete or can be safely converted to discrete. A classic example is balancing an inverted pendulum on a cart on a 2D plane. The angle of the pole can be represented as, instead of a continuous space $[0, 2\pi)$, a discrete space $[0, 2\pi,$ `step_number`]. This method has been proven to work with some trade-off of computational speed [3]. While it is possible to increase `step_number` to a large number and make the discrete space equivalent to continuous space in the computational world, the speed of the algorithm is traded off at an exponential rate.

### 2.1.2 Choose an action

With all action-state pairs listed out, the agent can interact with the environment. The process will start with the agent observing the current state of the environment; and based on its observation, an action is chosen to maximize the future reward. Initially, the agent would take action randomly; however, after each iteration and the Q-value gets updated (sec. 2.1.3), the agent will learn to take action more intentionally.

To avoid getting stuck in a local minimum, there has to be a method introduced to balance the exploration and exploitation rate of the agent. A simple method to do this is the epsilon-greedy algorithm [4]. In brief, epsilon $\epsilon$ is initialized as 1, and the agent would 100% choose action randomly (exploration). After each episode, $\epsilon$ will decay by a predetermined rate. At each episode, the agent would have $\epsilon$ chance to explore and $1 - \epsilon$ chance to exploit the environment.

In the program used for the project, epsilon is decayed following the following equation. $\epsilon_{min}$ is set up to ensure that the agent keeps exploring the environment at a small rate as time goes on.

$$\epsilon(t) = \epsilon_{min} + \frac{\epsilon_{max} - \epsilon_{min}}{\exp(t/\epsilon_{decay})}$$

### 2.1.3 Update Q-value

After the agent takes the action, the Q-value for that state-action pair needs to be calculated and saved for future iteration. The Bellman Equation is used for this step.

$$Q_{new}(s(t), a) = Q_{current}(s(t), a) + \alpha \left[ (R + \gamma * Q_{max}(s(t+1), a)) - Q_{current}(s(t), t) \right]$$

where:
- $Q(s(t), a) \rightarrow Q$ of state $s$ and action $a$

- $R \rightarrow$ reward when action $a$ is taken as state $s$
- $\gamma \rightarrow$ discount factor

The specific values for learning rate $\alpha$ and discount factor $\gamma$ are up to the user. However, there is always a trade-off when changing these values. With high $\alpha$, the agent would learn faster but may never get to the optimal result since the step size is too large.

## 2.2   Deep Q-Learning

Deep Q-learning is a variation of Q-learning. To combat a Q-Learning limitation of finite state and action space, the Q-value, instead of being calculated based on state-action pairs, is determined by a neural network.
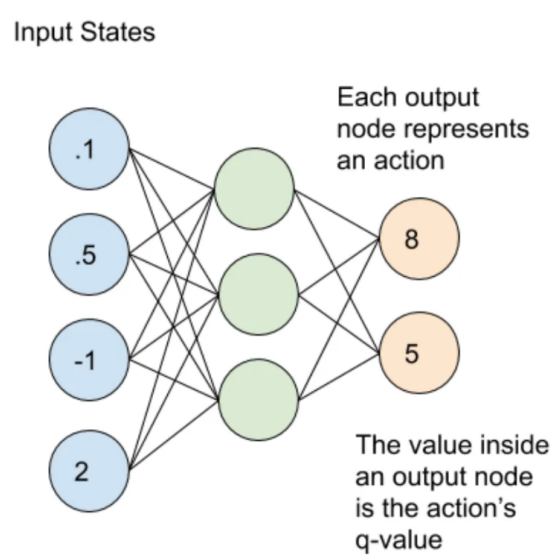


Figure 1: Deep Q Network maps states to actions and corresponding Q-values

To utilize Deep Q-Learning, instead of setting up a list of action-state pairs with randomly assigned Q-values, a neural network with random weights is set up. After each iteration, the weights of the neural network will be updated similar to Q-Learning. The complexity of the neural network can vary depending on the problem. For this project, a simple linear model consisting of two layers is implemented.

```
class DQN(nn.Module):
    def __init__(self, n_observations, n_actions):
        super(DQN, self).__init__()

        hidden_layer_1 = 128
        hidden_layer_2 = 128

        self.layer1 = nn.Linear(n_observations, hidden_layer_1)
        self.layer2 = nn.Linear(hidden_layer_1, hidden_layer_2)
        self.layer3 = nn.Linear(hidden_layer_2, n_actions)
```

```
12          # Called with either one element to determine next action, or a
    batch
13          # during optimization. Returns tensor([[left0exp,right0exp]...]).
14          def forward(self, x):
15              x = F.relu(self.layer1(x))
16              x = F.relu(self.layer2(x))
17              return self.layer3(x)
```

To add stability to the algorithm, a copy of the main work is created, called a target network. Instead of getting updated every episode like the main network, the target network is updated in a delayed manner. Specific rules to dictate when the target network is updated can be different depending on the program. For example, the target network can be updated after a specific number of episodes or it can be updated every episode but multiply with a scale factor. The program implemented for this project uses the second method.

```
1 # Soft update of the target network's weights
2 #                      + (1        )
3 target_net_state_dict = target_net.state_dict()
4 policy_net_state_dict = policy_net.state_dict()
5
6    for key in policy_net_state_dict:
7        target_net_state_dict[key] = policy_net_state_dict[key]*TAU +
    target_net_state_dict[key]*(1-TAU)
8        target_net.load_state_dict(target_net_state_dict)
```

The target network can also be removed entirely; and expectedly, that will support faster learning at the cost of stability. However, in my attempt, the network shows no positive learning after 1000 episodes.

# 3    Environment description

For this project, the `LunarLander-v2` from Gymnasium of Farama Foundation is implemented.

The environment tasks a space shuttle to land on a landing pad always at coordinates $(0,0)$. The terrain around the landing site is randomly generated. Initially, the spaceship is spawned at the top center of the viewport with a random force applied to its center of mass. To simplify the problem, there is no wind in the environment. This is accurate in the case of the moon but would need to change if the simulation is for another planet. The action space of the agent consists of four discrete actions:
- Do nothing
- Fire left orientation engine
- Fire main engine
- Fire right orientation engine.
The engine is assumed to be perfect and can only be turned fully on or off without variation.

In every frame, the observed state consists of six float values and two boolean values:
- $(x, y)$ coordinates of the spaceship
- $v(x, y)$ linear velocity of the spaceship
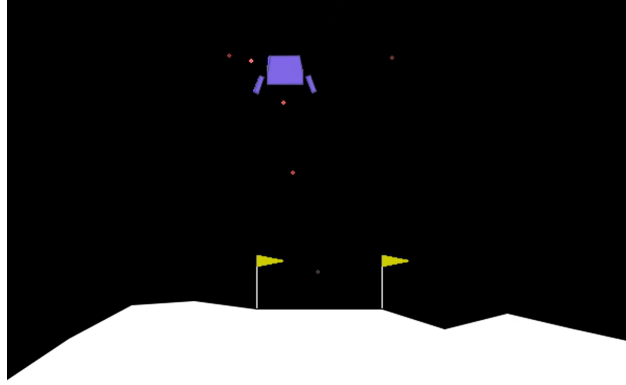- Angle of the spaceship

Figure 2: Example of an episode

- Angular velocity of the spaceship
- Whether the left leg is in contact with the ground (boolean).
- Whether the right leg is in contact with the ground (boolean)

After every step, a reward is granted. The total reward of an episode is the sum of the rewards for all the steps within that episode. For each step, the reward:
- is increased/decreased the closer/further the lander is to the landing pad.
- is increased/decreased the slower/faster the lander is moving.
- is decreased the more the lander is tilted (angle not horizontal).
- is increased by 10 points for each leg that is in contact with the ground.
- is decreased by 0.03 points each frame a side engine is firing.
- is decreased by 0.3 points each frame the main engine is firing.
The episode receives an additional reward of -100 or +100 points for crashing or landing safely respectively.

The episode ends when:
- the frame spent in the environment reaches 1000
- the spacecraft flies out of the environment window.

## 4    Results

Before analyzing the result, it is worthied to set expectations. The agent would not start training until its memory fills up a batch size, there should be little or no positive trends to the received rewards. Because there is punishment for each frame the engine is fired, the agent will take advantage of gravity to accelerate downward and only fire the engine last minute to correct itself. It should be noted that the landing gear is assumed to be indestructible, and there is no punishment for landing at a very high speed,

## 4.1 Episode rewards

The plot 3 shows that the agent successfully learns the environment. After 128 episodes, which filled up on batch, the agent starts to learn about the environment and quickly reaches 200 rewards, which is considered "solved" for this problem [5].
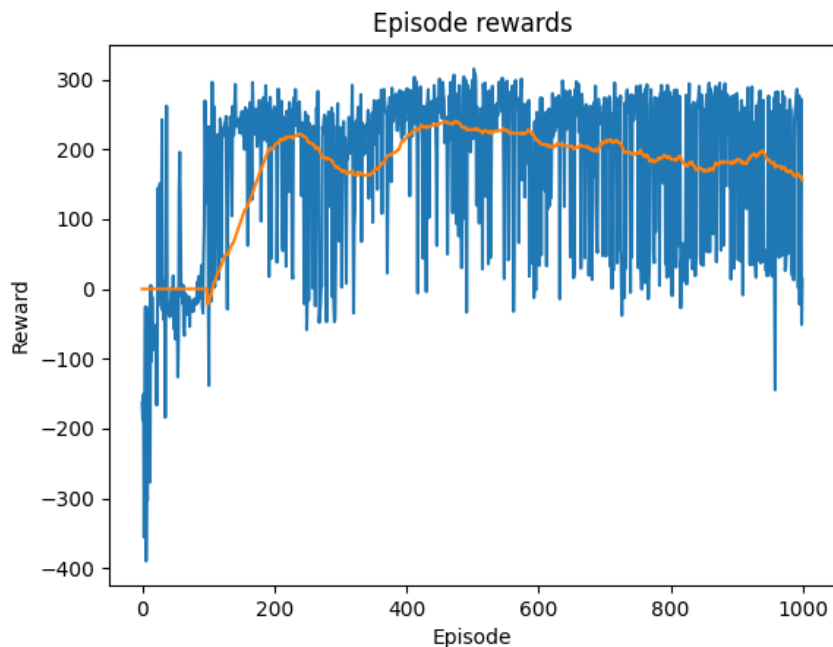


Figure 3: Rewards for 1000 episodes

After that, the reward curve went through a dip but quickly recovered later. This may be because the agent is trying to explore the environment instead of exploiting pre-existing strategies. However, because epsilon is relatively small at episode 300, it is unlikely that this would happen again if the program is executed again.

## 4.2 Episode duration

As mentioned above, there is a mild punishment for firing the engine, and there is no punishment for landing at a very high speed. The agent learned to land quicker as time went on (fig. 4). However, the rate of agents learning this is uncertain. In other attempts to train the neural network, the agent learns to cut down the time a lot slower. Usually, it takes about 400 episodes to see a clear drop in duration.
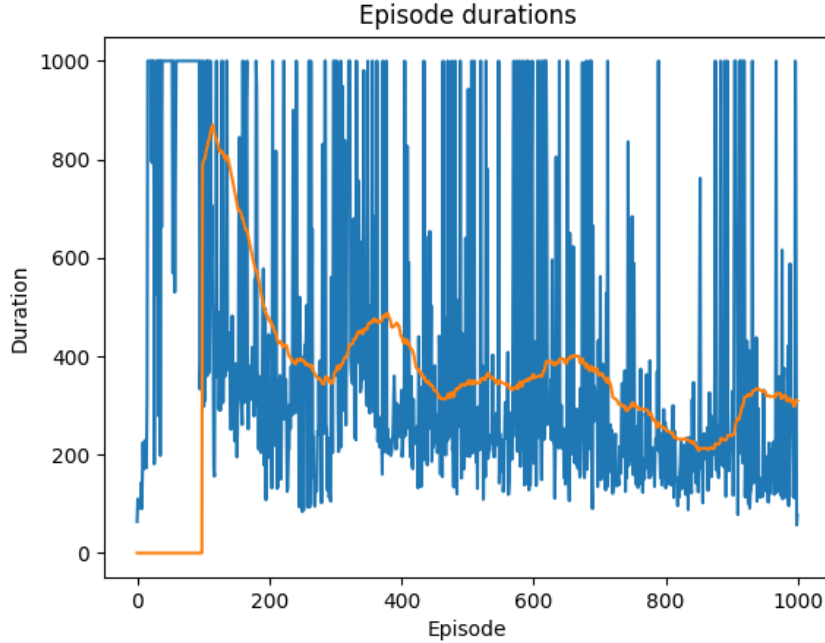
Figure 4: Durations of 1000 episodes

The speed at which the agent completes the task also comes with a trade-off. At about episode 400, the agent received the most reward. For this run to train the neural network, it seems like this is the optimal point where the agent landed the spaceship confidently and accurately. As time goes on, the episode durations keep decreasing (possibly a consequence of the rapid decrease in duration), and the episode rewards also experience a mild decrease. When the agent tries to land the spaceship too fast, its accuracy also decreases.

However, because of the discount factor $\gamma$, the calculated Q-value from experience would lose effect over time. The agent is expected to correct itself as time goes on.

# 5   Discussion

This project experimented with a recently developed reinforcement learning technique, Deep Q-Learning, to study an unknown environment. The implementation of the project shows promising results. The network solves the problem reasonably quickly and continues to improve even after an initial solution is found. It is possible to apply this technique to solve more complex systems even ones that are not fully understood.

For future study, some calibration to make the model more unstable may be introduced: removing the target network, using Q-Learning on continuous space, etc. These techniques are briefly mentioned in this report but were not implemented successfully in the program. A trial to apply the techniques to more complex simulations will also be considered. Simply introducing turbulence in the environment and control systems is an interesting path to

study the program more in-depth.

# Acknowledgement

# References

1. Watkins, C. J. C. H. & Dayan, P. Q-learning. *Machine Learning.* https://doi.org/ 10.1007/BF00992698 (1992).

2. Andersen, P. & Keiding, N. in *International Encyclopedia of the Social Behavioral Sciences* (eds Smelser, N. J. & Baltes, P. B.) 4946–4956 (Pergamon, Oxford, 2001). ISBN: 978-0-08-043076-8. https://www.sciencedirect.com/science/article/pii/ B0080430767021033.

3. Arts, L. *Comparing Discretization Methods for Applying Q learning in Continuous State-Action Space* 2017.

4. McCaffrey, J. D. *The Epsilon-Greedy Algorithm* https://jamesmccaffrey.wordpress. com/2017/11/30/the-epsilon-greedy-algorithm/.

5. Klimov, O. *Lunar Lander* https://gymnasium.farama.org/environments/box2d/ lunar_lander/#credits.

6. Paszke, A. & Towers, M. *Reinforcement Learning (DQN) Tutorial* https://pytorch. org/tutorials/intermediate/reinforcement_q_learning.html.

7. Foundation, F. *Gymnasium Documentation* https://gymnasium.farama.org/.

8. OpenAI. *ChatGPT* https://openai.com/chatgpt.

9. Google. *Gemini* https://gemini.google.com/.

# A  Project code

```python
1  import gymnasium as gym
2  import math
3  import random
4  import matplotlib
5  import matplotlib.pyplot as plt
6  from collections import namedtuple, deque
7  from itertools import count
8  import csv
9  # import pandas as pd
10
11 import torch
12 import torch.nn as nn
13 import torch.optim as optim
14 import torch.nn.functional as F
15
16 # env = gym.make("LunarLander-v2", render_mode = "human")
17 env = gym.make("LunarLander-v2")
18
19 # set up matplotlib
20 is_ipython = 'inline' in matplotlib.get_backend()
21 if is_ipython:
22     from IPython import display
23 plt.ion()
24
25 # use GPU if available
26 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
27 print("torch.cuda.get_device_name(0) = ", torch.cuda.get_device_name(0))
28
29 # %%
30
31 Transition = namedtuple("Transition", ("state", "action", "next_state", "
       reward"))
32
33 # %%
34
35 class ReplayMemory(object):
36     # create memory
37     def __init__(self, capacity):
38         self.memory = deque([], maxlen=capacity)
39
40     #
41     def push(self, *args):
42         """Save a transition"""
43         self.memory.append(Transition(*args))
44
45     def sample(self, batch_size):
46         # return random.sample(self.memory, batch_size)
47         return random.sample(self.memory, batch_size)
48
49     def __len__(self):
50         return len(self.memory)
```

11

```python
51
52  # %%
53
54  class DQN(nn.Module):
55
56      def __init__(self, n_observations, n_actions):
57          super(DQN, self).__init__()
58
59          hidden_layer_1 = 128
60          hidden_layer_2 = 128
61
62          self.layer1 = nn.Linear(n_observations, hidden_layer_1)
63          self.layer2 = nn.Linear(hidden_layer_1, hidden_layer_2)
64          self.layer3 = nn.Linear(hidden_layer_2, n_actions)
65
66      # Called with either one element to determine next action, or a batch
67      # during optimization. Returns tensor([[left0exp,right0exp]...]).
68      def forward(self, x):
69          x = F.relu(self.layer1(x))
70          x = F.relu(self.layer2(x))
71          return self.layer3(x)
72
73  # %%
74
75  # BATCH_SIZE is the number of transitions sampled from the replay buffer
76  # GAMMA is the discount factor as mentioned in the previous section
77  # EPS_START is the starting value of epsilon
78  # EPS_END is the final value of epsilon
79  # EPS_DECAY controls the rate of exponential decay of epsilon, higher
80      means a slower decay
80  # TAU is the update rate of the target network
81  # LR is the learning rate of the ``AdamW`` optimizer
82  BATCH_SIZE = 500
83  GAMMA = 0.99
84  EPS_START = 1.0
85  EPS_END = 0.05
86  EPS_DECAY = 1000
87  TAU = 0.005
88  LR = 1e-4
89
90  # Get number of actions from gym action space
91  n_actions = env.action_space.n
92  # Get the number of state observations
93  state, info = env.reset()
94  n_observations = len(state)
95
96  policy_net = DQN(n_observations, n_actions).to(device)
97  target_net = DQN(n_observations, n_actions).to(device)
98  target_net.load_state_dict(policy_net.state_dict())
99
100 optimizer = optim.AdamW(policy_net.parameters(), lr=LR, amsgrad=True)
101 memory = ReplayMemory(10000)
102
103
```

```python
104  steps_done = 0
105
106  # %%
107
108  def select_action(state):
109      global steps_done
110      sample = random.random()
111      eps_threshold = EPS_END + (EPS_START - EPS_END) * \
112          math.exp(-1. * steps_done / EPS_DECAY) # decaying epsilon
     exponentially
113      steps_done += 1
114      if sample > eps_threshold:
115          with torch.no_grad():
116              # t.max(1) will return the largest column value of each row.
117              # second column on max result is index of where max element
     was
118              # found, so we pick action with the larger expected reward.
119              return policy_net(state).max(1).indices.view(1, 1)
120      else:
121          return torch.tensor([[env.action_space.sample()]], device=device,
     dtype=torch.long)
122
123  # %%
124
125  episode_durations = []
126  episode_rewards = []
127
128  def plot_episode_data(episode_data, data_name:str, instant:int,
     show_result=False):
129      """
130      data_name should be capitalize
131      """
132      plt.figure(instant)
133      data_t = torch.tensor(episode_data, dtype=torch.float)
134      if show_result:
135          plt.title("Episode " + data_name)
136      else:
137          plt.clf()
138          plt.title('Training... (Episode', data_name + ")")
139      plt.xlabel('Episode')
140      plt.ylabel(data_name)
141      plt.plot(data_t.numpy())
142
143      # Take 100 episode averages and plot them too
144      if len(data_t) >= 100:
145          means = data_t.unfold(0, 100, 1).mean(1).view(-1)
146          means = torch.cat((torch.zeros(99), means))
147          plt.plot(means.numpy())
148
149      # pause a bit so that plots are updated
150      # can remove if num_episodes is too large
151      plt.pause(0.001)
152
153      # # use when writting code in ipynb
```

```python
154        # if is_ipython:
155        #     if not show_result:
156        #         display.display(plt.gcf())
157        #         display.clear_output(wait=True)
158        #     else:
159        #         display.display(plt.gcf())

161 # %%

163 # set up file to record training data
164 file_name = "training_data.csv"

166 def record_csv(episodes, rewards, durations, file_name=file_name):
167     # Open the file in write mode
168     with open(file_name, mode='w', newline='') as file:
169         writer = csv.writer(file)

171         # header
172         writer.writerow(['Episode', 'Reward', 'Duration'])

174         # record training data
175         for episode_val, reward_val, duration_val in zip(episodes, rewards
    , durations):
176             writer.writerow([episode_val.item(), reward_val.item(),
    duration_val])

178 # %%

180 def optimize_model():
181     # start training when memory is equal to BATCH_SIZE
182     if len(memory) < BATCH_SIZE:
183         return

185     transitions = memory.sample(BATCH_SIZE)
186     # Transpose the batch (see https://stackoverflow.com/a/19343/3343043
    for
187     # detailed explanation). This converts batch-array of Transitions
188     # to Transition of batch-arrays.
189     batch = Transition(*zip(*transitions))

191     # Compute a mask of non-final states and concatenate the batch
    elements
192     # (a final state would've been the one after which simulation ended)
193     non_final_mask = torch.tensor(tuple(map(lambda s: s is not None,
194                                         batch.next_state)), device=
    device, dtype=torch.bool)
195     non_final_next_states = torch.cat([s for s in batch.next_state
196                                         if s is not None])
197     state_batch = torch.cat(batch.state)
198     action_batch = torch.cat(batch.action)
199     reward_batch = torch.cat(batch.reward)

201     # Compute Q(s_t, a) - the model computes Q(s_t), then we select the
```

14

```python
202      # columns of actions taken. These are the actions which would've been
    taken
203      # for each batch state according to policy_net
204      state_action_values = policy_net(state_batch).gather(1, action_batch)
205
206      # Compute V(s_{t+1}) for all next states.
207      # Expected values of actions for non_final_next_states are computed
    based
208      # on the "older" target_net; selecting their best reward with max(1).
    values
209      # This is merged based on the mask, such that we'll have either the
    expected
210      # state value or 0 in case the state was final.
211      next_state_values = torch.zeros(BATCH_SIZE, device=device)
212      with torch.no_grad():
213          next_state_values[non_final_mask] = target_net(
    non_final_next_states).max(1).values
214      # Compute the expected Q values
215      expected_state_action_values = (next_state_values * GAMMA) +
    reward_batch
216
217      # Compute Huber loss
218      criterion = nn.SmoothL1Loss()
219      loss = criterion(state_action_values, expected_state_action_values.
    unsqueeze(1))
220
221      # Optimize the model
222      optimizer.zero_grad()
223      loss.backward()
224      # In-place gradient clipping
225      torch.nn.utils.clip_grad_value_(policy_net.parameters(), 100)
226      optimizer.step()
227
228  # %%
229
230  if torch.cuda.is_available():
231      num_episodes = 1000
232  else:
233      num_episodes = 50
234
235  for i_episode in range(num_episodes):
236      # Initialize the environment and get its state
237      state, info = env.reset()
238      state = torch.tensor(state, dtype=torch.float32, device=device).
    unsqueeze(0)
239
240      episode_reward = 0
241
242      for t in count():
243          action = select_action(state)
244          observation, reward, terminated, truncated, _ = env.step(action.
    item())
245          reward = torch.tensor([reward], device=device)
246          done = terminated or truncated
```

```python
247
248         if terminated:
249             next_state = None
250         else:
251             next_state = torch.tensor(observation, dtype=torch.float32,
     device=device).unsqueeze(0)
252
253         # Store the transition in memory
254         memory.push(state, action, next_state, reward)
255         episode_reward += reward
256
257         # Move to the next state
258         state = next_state
259
260         # Perform one step of the optimization (on the policy network)
261         optimize_model()
262
263         # Soft update of the target network's weights
264         #                    + (1        )
265         target_net_state_dict = target_net.state_dict()
266         policy_net_state_dict = policy_net.state_dict()
267         for key in policy_net_state_dict:
268             target_net_state_dict[key] = policy_net_state_dict[key]*TAU +
     target_net_state_dict[key]*(1-TAU)
269         target_net.load_state_dict(target_net_state_dict)
270
271         if done:
272             episode_durations.append(t + 1)
273             episode_rewards.append(episode_reward)
274             # plot_durations()
275             # plot_rewards()
276             print("Episode =", i_episode, "\tReward=%.2f" % episode_reward
     , "\tDuration=%.2f" % (t+1))
277             break
278
279 print('Complete')
280 # plot_durations(show_result=True)
281 # plot_rewards(show_result=True)
282 plot_episode_data(episode_rewards, "Reward", 1, True)
283 plot_episode_data(episode_durations, "Duration", 2, True)
284 record_csv(episodes=torch.arange(num_episodes), rewards=episode_rewards,
     durations=episode_durations)
285 plt.ioff()
286 plt.show()
```