

AI 3603 ARTIFICIAL INTELLIGENCE: PRINCIPLES AND TECHNIQUES

By: 陈路轩 (523030910014)

HW#: 2

2025 年 11 月 14 日

I. INTRODUCTION

A. Problem background

This assignment consists of three core tasks, designed to progressively build and analyze capabilities from basic tabular RL to Deep Reinforcement Learning (DRL):

- **Cliff-walking:** A classic grid-world problem. The agent must navigate a 12×4 grid from a start point to a goal point. The grid includes a "cliff" region, which provides a large negative reward and returns the agent to the start. The agent must learn to find an efficient path while avoiding the cliff.
- **Lunar Lander:** A more complex control problem. The agent must control a lander to touch down safely on a landing pad between two flags in a simulated lunar gravity environment. The state space is continuous (8 dimensions) and the action space is discrete (4 actions). This requires the agent to learn not just where to go, but how to control its thrusters for a smooth landing.

B. Task Objectives

This assignment has three core tasks, designed to progressively build a comprehensive path-planning system:

- **Task 1:** To implement the Sarsa, Q-Learning, and Dyna-Q algorithms. Then plot the episode reward curves, ϵ -decay curves, and visualize the final paths learned by the agents. Compare and analyze the performance of the three algorithms.
- **Task 2:** To read and understand the provided `dqn.py` code, adding comments to key sections. Then train and tune the DQN agent on the Lunar Lander environment, plotting the reward and ϵ -decay curves. Finally, visualize the trained agent's landing behavior.
- **Task 3:** Find and learn an exploration strategy other than ϵ -greedy.

II. ALGORITHM DESIGN AND IMPLEMENTATION

A. Task 1: Reinforcement Learning in Cliff-walking Environment

1. Description

Reinforcement Learning (RL) is a machine learning paradigm where an agent learns by trial and error through interactions with an environment. The goal is to acquire a policy to get as high as possible scores in the game. In this task, we implement RL agents based on Sarsa, Q-Learning, and Dyna-Q algorithms to find a safe path to the goal in a grid-shaped maze. The environment is a 12×4 grid map, and the agent is restricted to moving only upward, downward, leftward, and rightward.

2. Formulation

The core of these RL algorithms lies in their components: states, actions, and a reward function. The agent learns an action-value function, $Q(s, a)$, to determine the expected return of taking an action in a given state.

1. **State (s_t):** This value is an integer representing the agent's current coordinate (x, y) .
2. **Action (a_t):** This value is $a_t \in \{0, 1, 2, 3\}$, where the four integers represent the four moving directions (up, down, left, right) respectively.
3. **Reward (r):** This value represents the feedback from the environment. In this task, every step costs -1. Falling into the cliff gives a punishment of -100 and returns the agent to the starting point.

Q-value Update (Sarsa): The agent updates its Q-value estimation based on experience. The update rule for Sarsa (on-policy) is:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$$

Q-value Update (Q-Learning): The update rule for Q-Learning (off-policy) is:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

Model Learning and Planning (Dyna-Q): Dyna-Q integrates model-free interaction with model-based planning. In addition to the direct Q-Learning update from real experience (Direct RL), it involves:

1. **Model Learning:** The agent uses the experience (s, a, r, s') to update a model, $Model(s, a) \leftarrow (r, s')$.
2. **Planning:** The agent performs n planning steps. For each step, it randomly samples a previously observed state-action pair (s_{plan}, a_{plan}) , uses the model to get the simulated next state s'_{plan} and reward r_{plan} , and then applies the Q-Learning update rule:

$$Q(s_{plan}, a_{plan}) \leftarrow Q(s_{plan}, a_{plan}) + \alpha[r_{plan} + \gamma \max_{a'} Q(s'_{plan}, a') - Q(s_{plan}, a_{plan})]$$

In these formulas, s is the current state, a is the current action, r is the reward received, s' is the next state, and a' is the next action. α is the learning rate, and γ is the reward decay.

3. Implementation

a. Core Data Structures:

- (a) **Agent Class:** Separate classes are defined for each algorithm (SarsaAgent, QLearningAgent, Dyna_QAgent) in `agent.py`. Each instance stores hyperparameters like `alpha`, `gamma`, and `epsilon`.
- (b) **Q-Table:** A data structure (typically a 2D numpy array) within each agent, mapping state-action pairs to their estimated Q-values. It is initialized with the state dimension and number of actions.
- (c) **Model (only for Dyna-Q):** A data structure specific to `Dyna_QAgent` used to store experiences (s, a, r, s') . This model is used for planning by simulating interactions.

b. Algorithm Execution Flow:

- (a) **Initialization:** Create the Gym environment (`gym.make()`). Then, construct the agent with its hyperparameters.
- (b) **Main Training Loop:** The agent is trained over a fixed number of episodes (e.g., 1000).
- (c) **Episode Reset:** At the start of each episode, the environment is reset (`env.reset()`) to get the initial state `s`.
- (d) **Step Interaction Loop:** Within an episode, the agent interacts with the environment for a maximum number of steps or until the terminal state is reached (`isdone`).
- (e) **Action Selection:** The agent chooses an action `a` based on state `s` using an ϵ -greedy policy (Explore randomly with probability `epsilon`; otherwise, exploit the currently known best action.) (`agent.choose_action(s)`).
- (f) **Environment Step:** The chosen action `a` is sent to the environment (`env.step(a)`), which returns the next state `s_`, reward `r`, and done flag `isdone`.
- (g) **Learning (Q-Update):** The agent's `learn` function is called with the experience tuple.
 - **Sarsa:** Requires the next action `a_` chosen by the policy at `s_`. The update is `agent.learn(s, a, r, s_, a_)`.
 - **Q-Learning / Dyna-Q:** The update is based on the max Q-value at `s_`. The call is `agent.learn(s, a, r, s_)`. Dyna-Q also performs n planning steps internally during this call.
- (h) **State Transition:** The current state `s` is updated to `s_`. (For Sarsa, `a` is also updated to `a_`).
- (i) **Epsilon Decay:** After each episode, the `epsilon` value is reduced via `agent.decay_epsilon()` to shift from exploration to exploitation.
- (j) **Result Visualization:** After training, the epsilon decay curve and moving average of rewards are plotted. A final test run with $\epsilon = 0$ is performed, and the resulting path is visualized.

4. Parameter Tuning

After multiple rounds of experimentation and parameter tuning, the current set of values (**alpha** = 0.1, **gamma** = 0.9, **epsilon** = 1.0, **epsilon_decay** = 0.99, **epsilon_min** = 0.01, and **n_planning_steps** = 100) was determined to be a good solution for the Cliff-walking task. Training was conducted over 1000 episodes.

During the tuning process, we observe the impact of each hyperparameter. The **alpha** (learning rate) parameter dictates the convergence speed. An excessively large alpha caused Q-value estimates to oscillate and fail to converge, while a value that was too small resulted in slow learning. The **gamma** (discount factor) influences the agent’s foresight. A value close to 1, such as 0.9, is crucial for this task, as it encourages the agent to value long-term rewards (reaching the exit) over short-term costs (the -1 for each step).

The **epsilon** decay schema manages the critical balance between exploration and exploitation. An initial **epsilon** of 1.0 ensures full exploration at the beginning of training. The **epsilon_decay** rate of 0.99 ensures a gradual transition, allowing the agent to exploit known good paths as training progresses. For Dyna-Q, the **n_planning_steps** parameter directly impacts training efficiency. A higher value allows the agent to learn more from each real experience, leading to much faster convergence in terms of episodes, at the cost of more computation per step.

The selected parameters achieve an effective balance between sufficient exploration to discover the optimal path, stable Q-value convergence, and an efficient learning process across all three implemented algorithms.

5. Result

Here are the Epsilon decay curves, average reward curves and final paths for the three algorithms:

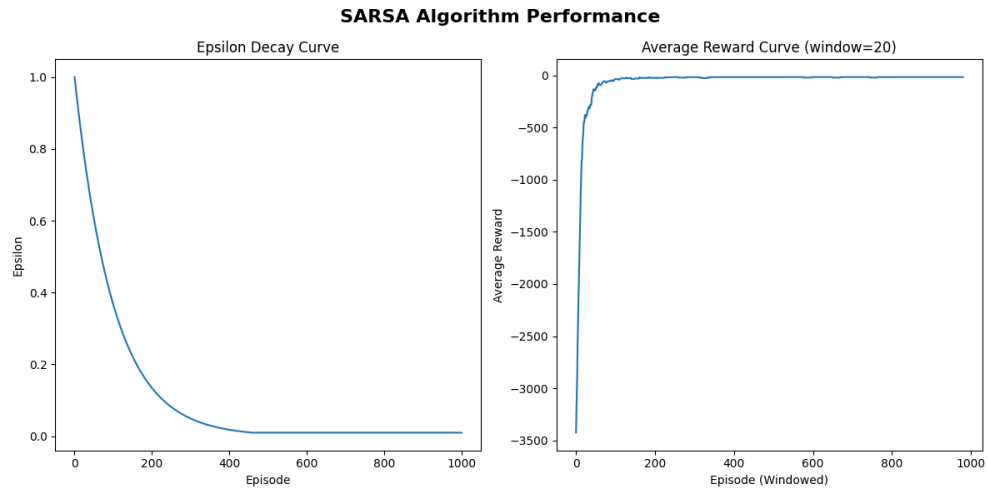


Figure 1: SARSA Performance



Figure 2: SARSA Path

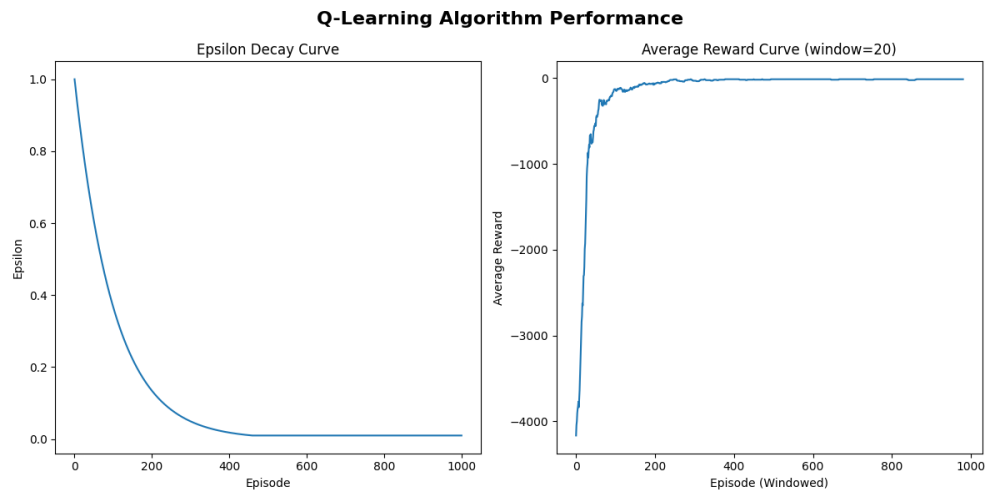


Figure 3: Qlearning Performance



Figure 4: Qlearning Path

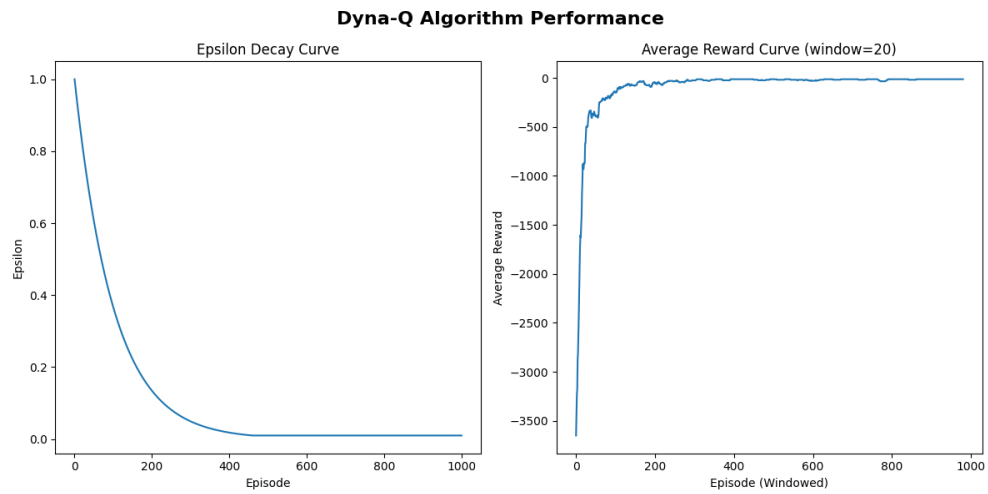


Figure 5: Dyna-Q Performance



Figure 6: Dyna-Q Path

1. Path difference between Sarsa and Q-learning:

The paths generated by Sarsa and Q-learning exhibit notable differences due to their underlying learning strategies. Q-Learning finds the optimal path, which runs directly along the edge of the cliff with a reward of -13. While SARSA finds a "safer" and longer path that detours to the topmost row of the grid, far away from the cliff with a cost of -17.

This is because Q-learning is an off-policy algorithm that learns the optimal policy regardless of the agent's actions, leading it to exploit the cliff-edge path for maximum reward. In contrast, Sarsa is an on-policy algorithm that learns the value of the policy being followed, which includes the risk of falling into the cliff. As a result, Sarsa tends to favor safer paths that avoid high-risk areas, even if they are longer.

- SARSA(on-policy): Its update must account for the actual next action (a') chosen by its ϵ -greedy policy⁵. When on the cliff's edge, the ϵ -greedy policy has a non-zero chance of randomly selecting "down," incurring a massive -100 penalty. SARSA learns to factor in this exploration risk. It converges on a more "conservative" policy, preferring the longer path (more -1 penalties) to avoid the risk of being near the cliff.
- Q-Learning(off-policy): Its update is based on the theoretical best-possible next action ($\max a'$). It learns the value of the optimal policy without regard for the risks taken during exploration. Therefore, it finds the shortest path, ignoring the risk that the ϵ -greedy exploration policy might accidentally step off the cliff.

2. Training efficiency between model-based RL (dyna-Q) and model-free algorithms (Sarsa or Q-learning)

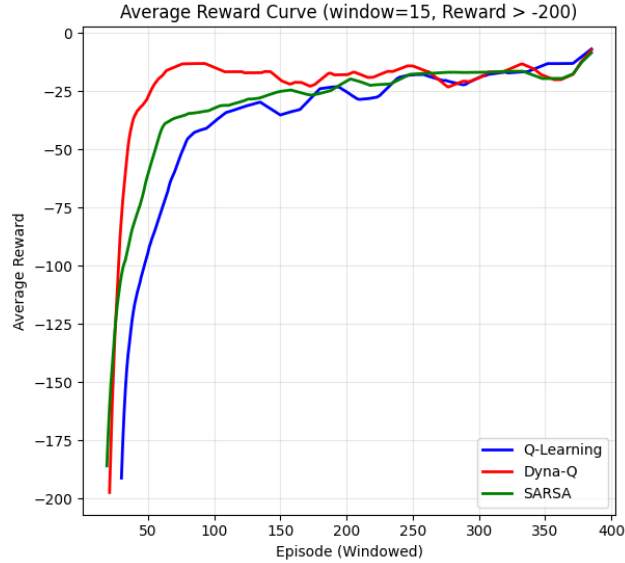


Figure 7: training efficiency comparison

As shown in the graph, Dyna-Q (red line) is significantly more training-efficient than both Q-Learning (blue line) and SARSA (green line). Dyna-Q's average reward curve rises the fastest, reaching a high-reward plateau within approximately 50-75 episodes. In contrast, SARSA and Q-Learning require 100-150 episodes to reach a similar level of performance.

This difference is mainly because Dyna-Q is a model-based algorithm, while SARSA and Q-Learning are model-free.

- Model-Free (Sarsa/Q-Learning): These agents get only one Q-value update for each step of real interaction with the environment, based on the single experience tuple (s, a, r, s') .
- Model-Based (Dyna-Q): Dyna-Q simultaneously learns a model of the environment (i.e., $Model(s, a) \rightarrow r, s'$). For each single step of real experience, the Dyna-Q agent performs one direct Q-update (like Q-Learning) and additionally performs n planning steps (where $n = 100$ in the code). In these 100 planning steps, it uses its learned model to simulate experiences and updates its Q-table with these simulated experiences.

Therefore, Dyna-Q gets 101 learning opportunities per real interaction, while SARSA and Q-Learning gets only one. This allows Dyna-Q to squeeze much more learning out of each piece of real experience, leading to greater efficiency in terms of episodes (interactions).

B. Task 2: Deep Reinforcement Learning

1. Description

In this task, we implement a Deep Q-Network (DQN) agent to solve a more complex control problem: the "LunarLander-v2" gym environment. The goal is to control a spaceship and land it smoothly between two flags on the moon's surface. Unlike the previous task, the state space in this environment is continuous and high-dimensional (an 8-dimensional vector), which makes a tabular Q-table infeasible. Therefore, we use a neural network as a function approximator to estimate the action-value function $Q(s, a)$. The provided code, `dqn.py`, implements a DQN agent with enhancements such as a target network and experience replay.

2. Formulation

The core of the DQN algorithm is using a deep neural network to learn the optimal Q-value function.

1. **State (s_t):** This value is an 8-dimensional continuous vector representing the lander's physical status:

$$s_t = [x, y, v_x, v_y, \theta, \omega, leg_1, leg_2]$$

where (x, y) are coordinates, (v_x, v_y) are linear velocities, θ is the angle, ω is the angular velocity, and (leg_1, leg_2) are booleans indicating ground contact for each leg.

2. **Action (a_t):** This is an integer $a_t \in \{0, 1, 2, 3\}$, representing four discrete actions: do nothing, fire left orientation engine, fire main engine, and fire right orientation engine.
3. **Q-Value Approximation (Q-Network)** The action-value function $Q(s, a; \theta)$ is approximated by a neural network with parameters θ . The network architecture implemented in `dqn.py` is a Multi-Layer Perceptron (MLP):

- **Input Layer:** 8 neurons (matching the state dimension)
- **Hidden Layer 1:** 120 neurons with ReLU activation
- **Hidden Layer 2:** 84 neurons with ReLU activation
- **Output Layer:** 4 neurons (matching the action dimension), providing the Q-value for each possible action from the input state

4. **Loss Function and Update:** The network is trained by sampling a mini-batch of experiences (s, a, r, s', d) from a replay Buffer. The network's parameters θ are updated by minimizing the Mean Squared Error (MSE) loss:

$$L(\theta) = \mathbb{E}_{(s, a, r, s', d) \sim \text{Buffer}} [(y - Q(s, a; \theta))^2]$$

where y is the TD target. Unlike standard DQN, the implementation in `dqn.py` (after we modified it) uses the **Double DQN** rule to calculate the target, which helps reduce overestimation of Q-values. Double DQN decouples action selection from action evaluation:

- (a) The **main network** ($Q(s, a; \theta)$) is used to **select** the best action a^* in the next state s' :

$$a^* = \arg \max_{a'} Q(s', a'; \theta)$$

- (b) The **target network** ($Q(s, a; \theta')$) is used to **evaluate** the value of that action a^* .

The TD target y is therefore calculated as:

$$y = r + \gamma Q(s', a^*; \theta') = r + \gamma Q(s', \arg \max_{a'} Q(s', a'; \theta); \theta')$$

3. Implementation

a. Core Data Structures:

- (a) **QNetwork Class:** A neural network class is defined using `torch.nn`. Each instance is a Multi-Layer Perceptron (MLP) with an 8-neuron input layer, two hidden layers (120 and 84 neurons) with ReLU activation, and a 4-neuron output layer.
- (b) **Target Network:** A second instance of `QNetwork` is created, named `target_network`. Its weights are periodically synchronized with the main `q_network` (every `target_network_frequency` steps) to stabilize training.
- (c) **Replay Buffer:** An instance of `ReplayBuffer` from the `stable_baselines3` library is implemented. It stores up to `buffer_size` (100,000) of past transitions (state, action, reward, next state, done flag).
- (d) **Optimizer:** A `torch.optim.Adam` optimizer is initialized to update the parameters of the main `q_network`.

b. Algorithm Execution Flow:

- (a) **Initialization:** Parse command-line arguments, set random seeds for reproducibility, create the Gym environment using `make_env`, and initialize the `q_network`, `target_network`, `optimizer`, and `replay_buffer`. A `TensorBoard SummaryWriter` is also set up for logging.
- (b) **Main Loop:** The agent interacts with the environment in a single loop iterating for `total_timesteps`.
- (c) **Epsilon Calculation:** At each global step, the exploration probability `epsilon` is calculated using a `linear_schedule` function. It decays linearly from `start_e` to `end_e` over the first `exploration_fraction` of total timesteps.
- (d) **Action Selection:** An ϵ -greedy policy is used. A random number is compared to `epsilon`. If less, a random action is sampled (exploration). Otherwise, the main `q_network` predicts Q-values for the current state, and the action with the maximum Q-value is chosen (exploitation).
- (e) **Environment Interaction:** The chosen action is passed to `envs.step()`, which returns the `next_obs`, `rewards`, `done` flag, and `infos` dictionary.

- (f) **Store Experience:** The transition tuple (`obs`, `next_obs`, `actions`, `rewards`, `done`) is added to the `replay_buffer`.
- (g) **State Transition:** The current observation `obs` is updated to `next_obs`. If the episode is `done`, the environment is automatically reset.
- (h) **Learning Trigger:** The learning process (Q-update) begins only after `global_step` exceeds `learning_starts` and executes every `train_frequency` steps.
- (i) **Sample Batch:** A mini-batch of `batch_size` transitions is randomly sampled from the `replay_buffer`.
- (j) **TD Target Calculation:** The TD target is calculated using the Double DQN rule. First, the main `q_network` is used to select the best action for the next state (`next_actions`). Then, the `target_network` is used to evaluate the Q-value of that selected action (`target_max`). The final target is computed as $y = r + \gamma \times \text{target_max} \times (1 - d)$.
- (k) **Loss Calculation:** The `F.mse_loss` (Mean Squared Error) is computed between the `td_target` (y) and the `old_val` (the Q-value of the action actually taken, predicted by the main `q_network`).
- (l) **Optimization:** The optimizer’s gradients are cleared (`zero_grad()`), the loss is backpropagated (`loss.backward()`), gradient clipping is applied using `clip_grad_norm_` to prevent exploding gradients, and the optimizer updates the `q_network` weights (`optimizer.step()`).
- (m) **Target Network Update:** Every `target_network_frequency` steps, the `target_network`’s weights are synchronized by loading the state dictionary from the main `q_network`.

4. Parameter Tuning

After multiple rounds of experimentation and parameter tuning, the current set of values (`learning_rate` = 2.5e-4, `gamma` = 0.995, `buffer_size` = 100,000, `batch_size` = 256, and `target_network_frequency` = 500) was determined to be a near-optimal solution for the LunarLander-v2 task. The agent was trained for a total of 500,000 timesteps.

A high gamma value encourages the agent to prioritize long-term rewards (achieving a successful landing) over immediate, small rewards. The selected `batch_size` (256) offered a good balance between stability and computational speed. The `buffer_size` was set to 100,000 to ensure a large and diverse set of experiences.

The `epsilon` decay schema was tuned to start at 1.0 (full exploration) and decay linearly to 0.01 over 15% of the total timesteps. This extended exploration phase allows the agent to discover variable landing strategies before exploiting them. Finally, a `max_grad_norm` of 10.0 was applied to clip gradients, preventing gradient explosion and stabilizing the learning process.

The selected parameters achieve an excellent balance between stable learning, sufficient exploration, and efficient convergence, leading to a high-performance agent.

5. Result

We use the average reward over 200 test episodes as the final performance metric.

In our experiments, by only tuning hyperparameters on the original `dqn.py` code (such as `gamma`, `epsilon decay`, etc.), the best average reward we achieved was approximately 220. Subsequently, by tuning parameters and optimizing the network structure, the average reward increased to around 254. Finally, building on the previous optimizations, we implemented key techniques such as Double DQN and gradient clipping. This approach yielded the best performance, reaching an average reward of approximately 276, and it is the solution we ultimately adopted. Additionally, we also attempted to use multi-frame state stacking instead of a single frame for training, but the results were not satisfactory.

We have placed a video demonstration in the "videos" folder, which achieved a score of 281.45.

Note: You may find two versions of `dqn.py` in the submission: `dqn_original.py` (the original code provided) and `dqn.py` (our final modified version). That's because the assignment requires to make comments on the original one, so I wrote comments for both of them. The results mentioned above are based on our modified version.

C. Task 3: Improve Exploration Schema

Beyond the ϵ -greedy strategy, **Upper Confidence Bound** (UCB) is a more efficient and directed exploration method.

1. Idea

The core idea of UCB is "Optimism in the face of uncertainty". Unlike ϵ -greedy, which performs completely random exploration, UCB explores based on the degree of uncertainty about the value of each action.

ϵ -greedy might randomly select an action it already knows is bad. UCB, however, selects the action that it is most uncertain about but has the potential to be the best.

The UCB algorithm balances "Exploitation" and "Exploration" using a formula. When selecting an action a in state s , it chooses the action that maximizes the following expression:

$$a_t = \arg \max_a \left[Q(s, a) + C \sqrt{\frac{\ln t}{N(s, a)}} \right]$$

1. $Q(s, a)$ (**Exploitation Term**): This is the agent's current estimated value of taking action a in state s . The higher this value, the more the agent wants to exploit this action.
2. $C \sqrt{\frac{\ln t}{N(s, a)}}$ (**Exploration Term**): This is the uncertainty bonus.
 - t is the total number of decision steps (or episodes) so far.
 - $N(s, a)$ is the number of times action a has been selected in state s .
 - C is a constant that balances the weight of exploitation and exploration.

How it works:

- If $N(s, a)$ is small (i.e., this action has been tried infrequently in this state), the denominator is small, causing the exploration term to become very large. This strongly encourages the agent to try this action it knows little about.
- If $N(s, a)$ is large (i.e., this action has been tried many times), the denominator is large, and the exploration term approaches 0. The choice will then be based primarily on the actual $Q(s, a)$ value.
- As the total steps t increases, $\ln t$ grows slowly, ensuring that the agent never completely stops exploring.

2. Pros

1. **Efficient and Directed Exploration:** UCB's exploration is not random but strategic. It prioritizes exploring actions with the highest "potential" (highest uncertainty) rather than wasting time on actions known to be suboptimal.

2. **No Epsilon Decay Tuning:** The effectiveness of ϵ -greedy heavily relies on the initial ϵ value and a complex decay schedule. UCB adjusts its exploration automatically based on visit counts, and the parameter C is relatively less sensitive.

3. Cons

1. **Difficult to Scale to Large State Spaces:** The form of UCB requires maintaining an exact visit count $N(s, a)$ for every state-action pair (s, a) . This is feasible in tabular environments (like Task 1, Cliff-walking) but is impossible in environments with large or continuous state spaces (like Task 2, LunarLander-v2).
2. **Complex to Combine with Deep Learning:** In DQN, states are high-dimensional vectors, and a neural network cannot directly store $N(s, a)$. For example, s is not an integer, but rather an 8-dimensional vector of floating-point numbers (e.g., $[0.123, -0.456, 0.789, \dots, 1.0]$). In two consecutive decisions, it is virtually impossible for the agent to visit the exact same floating-point state. To apply UCB's "optimism" in Deep RL, researchers must use very complex methods to estimate uncertainty, such as:
 - **Pseudo-Counts:** Using the novelty of states to estimate $N(s, a)$.
 - **Bayesian Neural Networks:** Using NNs to output a distribution over Q-values rather than a single point estimate, thereby quantifying uncertainty.

These methods are all significantly more complex to implement than the ϵ -greedy strategy.

III. DISCUSSION AND CONCLUSION

A. Discussion of Findings

In **Task 1 (Cliff-walking)**, we observed the fundamental differences between RL algorithms. The Sarsa (on-policy) agent learned a "safer" path (Reward: -17), while the Q-Learning (off-policy) agent found the optimal but riskier path along the cliff edge (Reward: -13). Furthermore, Dyna-Q (model-based) showed superior sample efficiency, converging in far fewer episodes than the model-free methods, as it used its learned model for 100 planning steps per real interaction.

In **Task 2 (LunarLander-v2)**, the high-dimensional continuous state space required a DQN. We tuned the agent, using the average reward over 200 test episodes as our metric. Our final solution, which achieved an average reward of 276, was built by applying both Double DQN and gradient clipping.

In **Task 3**, we analyzed UCB as an alternative exploration strategy. UCB relies on precise visit counts $N(s, a)$, which is infeasible when the state is a high-dimensional float vector where the exact same state is almost never revisited. This makes the simplicity of ϵ -greedy far more practical for DQN.

B. Conclusion

This report successfully demonstrated the trade-offs between on-policy (safer) and off-policy (more optimal) methods, as well as the sample efficiency gains of model-based (Dyna-Q) over model-free (Sarsa/Q-Learning) algorithms. For the complex Lunar Lander task, we found that modern enhancements like Double DQN and gradient clipping are essential for achieving stable, high-performance results.

IV. APPENDIX

附录 A: Source Code for SARSA:

```
1  # -*- coding:utf-8 -*-
2  # Train Sarsa in cliff-walking environment
3  import math, os, time, sys
4  import numpy as np
5  import random
6  import gym
7  from agent import SarsaAgent
8
9  # construct the environment
10 env = gym.make("CliffWalking-v0")
11 # get the size of action space
12 num_actions = env.action_space.n
13 all_actions = np.arange(num_actions)
14 # set random seed and make the result reproducible
15 RANDOM_SEED = 0
16 env.seed(RANDOM_SEED)
17 random.seed(RANDOM_SEED)
18 np.random.seed(RANDOM_SEED)
19
20 ##### START CODING HERE #####
21 num_states = env.observation_space.n
22
23 # construct the intelligent agent.
24 agent = SarsaAgent(
25     all_actions=all_actions,
26     state_dim=num_states,
27     alpha=0.1,
28     gamma=0.9,
29     epsilon=1.0,
30     epsilon_decay=0.99,
31     epsilon_min=0.01,
32 )
33
34 episode_rewards = []
35 epsilon_values = []
36
37 # start training
38 for episode in range(1000):
39     # record the reward in an episode
40     episode_reward = 0
41     # reset env
```

```

42     s = env.reset()
43     # agent interacts with the environment
44     a = agent.choose_action(s)
45     for iter in range(500):
46         s_, r, isdone, info = env.step(a)
47         a_ = agent.choose_action(s_)
48         agent.learn(s, a, r, s_, a_)
49         s = s_
50         a = a_
51         episode_reward += r
52         if isdone:
53             # time.sleep(0.1)
54             break
55         episode_rewards.append(episode_reward)
56         epsilon_values.append(agent.epsilon)
57         agent.decay_epsilon()
58
59     if (episode + 1) % 50 == 0:
60         print(
61             "episode:",
62             episode + 1,
63             "episode_reward:",
64             episode_reward,
65             "epsilon:",
66             agent.epsilon,
67         )
68
69     print("\ntraining over\n")
70
71     # close the render window after training.
72     env.close()
73
74     import matplotlib.pyplot as plt
75     plt.figure(figsize=(12, 6))
76     plt.suptitle('SARSA Algorithm Performance', fontsize=16, fontweight='bold')
77     plt.subplot(1, 2, 1)
78     plt.plot(epsilon_values)
79     plt.title('Epsilon Decay Curve')
80     plt.xlabel('Episode')
81     plt.ylabel('Epsilon')
82
83     def moving_average(data, window_size=20):
84         return np.convolve(data, np.ones(window_size)/window_size, mode='valid')
85
86     avg_rewards = moving_average(episode_rewards)

```

```

87     plt.subplot(1, 2, 2)
88     plt.plot(avg_rewards)
89     plt.title(f'Average_Reward_Curve_(window={20})')
90     plt.xlabel('Episode_(Windowed)')
91     plt.ylabel('Average_Reward')
92
93     plt.tight_layout()
94     plt.show()
95
96     s = env.reset()
97     agent.epsilon = 0.0
98     episode_reward = 0
99     isdone = False
100    path = [s]
101
102    while not isdone:
103        env.render()
104        time.sleep(0.3)
105        a = agent.choose_action(s)
106        s, r, isdone, info = env.step(a)
107        path.append(s)
108        episode_reward += r
109
110    print(f"Test_complete!_Final_path_reward:_{episode_reward}")
111
112    img = env.render(mode='rgb_array')
113    env.close()
114    plt.figure(figsize=(12, 4))
115    plt.imshow(img)
116    coords = [(divmod(s, 12)[1] * img.shape[1] / 12 + img.shape[1] / 24,
117               divmod(s, 12)[0] * img.shape[0] / 4 + img.shape[0] / 8) for s in path]
118    plt.plot([c[0] for c in coords], [c[1] for c in coords], 'r-', linewidth=3, marker
              ='o', markersize=5)
119    plt.title(f'SARSA_Path_(Reward:_{episode_reward})', fontsize=14, fontweight='bold'
              )
120    plt.axis('off')
121    plt.savefig('sarsa_path.png', dpi=150, bbox_inches='tight')
122    plt.show()
123    ##### END CODING HERE #####

```

Listing 1: Source Code for SARSA (cliff_walk_sarsa.py)

附录 B: Source Code for Q-Learning:

```
1  # -*- coding:utf-8 -*-
2  # Train Q-Learning in cliff-walking environment
3  import math, os, time, sys
4  import numpy as np
5  import random
6  import gym
7  from agent import QLearningAgent
8  ##### START CODING HERE #####
9  # This code block is optional. You can import other libraries or define your
   utility functions if necessary.
10
11 ##### END CODING HERE #####
12
13 # construct the environment
14 env = gym.make("CliffWalking-v0")
15 # get the size of action space
16 num_actions = env.action_space.n
17 all_actions = np.arange(num_actions)
18 # set random seed and make the result reproducible
19 RANDOM_SEED = 0
20 env.seed(RANDOM_SEED)
21 random.seed(RANDOM_SEED)
22 np.random.seed(RANDOM_SEED)
23
24 ##### START CODING HERE #####
25 num_states = env.observation_space.n
26 # construct the intelligent agent.
27 agent = QLearningAgent(
28     all_actions=all_actions,
29     state_dim=num_states,
30     alpha=0.1,
31     gamma=0.9,
32     epsilon=1.0,
33     epsilon_decay=0.99,
34     epsilon_min=0.01,
35 )
36
37 episode_rewards = []
38 epsilon_values = []
39
40 # start training
41 for episode in range(1000):
42     # record the reward in an episode
```

```

43     episode_reward = 0
44     # reset env
45     s = env.reset()
46
47     # agent interacts with the environment
48     for iter in range(500):
49         # choose an action
50         a = agent.choose_action(s)
51         s_, r, isdone, info = env.step(a)
52         agent.learn(s, a, r, s_)
53         s = s_
54         # update the episode reward
55         episode_reward += r
56         if isdone:
57             time.sleep(0.1)
58             break
59         episode_rewards.append(episode_reward)
60         epsilon_values.append(agent.epsilon)
61         agent.decay_epsilon()
62
63     if (episode + 1) % 50 == 0:
64         print(
65             "episode:",
66             episode + 1,
67             "episode_reward:",
68             episode_reward,
69             "epsilon:",
70             agent.epsilon,
71         )
72     print('\ntraining over\n')
73
74     # close the render window after training.
75     env.close()
76
77
78     import matplotlib.pyplot as plt
79     plt.figure(figsize=(12, 6))
80     plt.suptitle('Q-Learning Algorithm Performance', fontsize=16, fontweight='bold')
81     plt.subplot(1, 2, 1)
82     plt.plot(epsilon_values)
83     plt.title('Epsilon Decay Curve')
84     plt.xlabel('Episode')
85     plt.ylabel('Epsilon')
86
87     def moving_average(data, window_size=20):

```

```

88     return np.convolve(data, np.ones(window_size)/window_size, mode='valid')
89
90     avg_rewards = moving_average(episode_rewards)
91     plt.subplot(1, 2, 2)
92     plt.plot(avg_rewards)
93     plt.title(f'Average_Reward_Curve_(window={20})')
94     plt.xlabel('Episode_(Windowed)')
95     plt.ylabel('Average_Reward')
96
97     plt.tight_layout()
98     plt.show()
99
100    s = env.reset()
101    agent.epsilon = 0.0
102    episode_reward = 0
103    isdone = False
104    path = [s]
105
106    while not isdone:
107        env.render()
108        time.sleep(0.3)
109        a = agent.choose_action(s)
110        s, r, isdone, info = env.step(a)
111        path.append(s)
112        episode_reward += r
113
114    print(f"Test_complete!_Final_path_reward:_{episode_reward}")
115
116    img = env.render(mode='rgb_array')
117    env.close()
118    plt.figure(figsize=(12, 4))
119    plt.imshow(img)
120    coords = [(divmod(s, 12)[1] * img.shape[1] / 12 + img.shape[1] / 24,
121               divmod(s, 12)[0] * img.shape[0] / 4 + img.shape[0] / 8) for s in path]
122    plt.plot([c[0] for c in coords], [c[1] for c in coords], 'r-', linewidth=3, marker
              = 'o', markersize=5)
123    plt.title(f'Q-Learning_Path_(Reward:_{episode_reward})', fontsize=14, fontweight='
              bold')
124    plt.axis('off')
125    plt.savefig('qlearning_path.png', dpi=150, bbox_inches='tight')
126    plt.show()
127    ##### END CODING HERE #####

```

Listing 2: Source Code for Q-learning (cliff_walk_qlearning.py)

附录 C: Source Code for Dyna-Q:

```
1      # -*- coding:utf-8 -*-
2      # Train Q-Learning in cliff-walking environment
3      import math, os, time, sys
4      import numpy as np
5      import random
6      import gym
7      from agent import Dyna_QAgent
8
9      # construct the environment
10     env = gym.make("CliffWalking-v0")
11     # get the size of action space
12     num_actions = env.action_space.n
13     all_actions = np.arange(num_actions)
14     # set random seed and make the result reproducible
15     RANDOM_SEED = 0
16     env.seed(RANDOM_SEED)
17     random.seed(RANDOM_SEED)
18     np.random.seed(RANDOM_SEED)
19
20     ##### START CODING HERE #####
21     num_states = env.observation_space.n
22
23     # construct the intelligent agent.
24     agent = Dyna_QAgent(
25         all_actions=all_actions,
26         state_dim=num_states,
27         alpha=0.1,
28         gamma=0.9,
29         epsilon=1.0,
30         epsilon_decay=0.99,
31         epsilon_min=0.01,
32         n_planning_steps=100,
33     )
34
35     episode_rewards = []
36     epsilon_values = []
37
38     # start training
39     for episode in range(1000):
40         # record the reward in an episode
41         episode_reward = 0
42         # reset env
43         s = env.reset()
```



```

44
45     # agent interacts with the environment
46     for iter in range(500):
47         # choose an action
48         a = agent.choose_action(s)
49         s_, r, isdone, info = env.step(a)
50         agent.learn(s, a, r, s_)
51         s = s_
52         # update the episode reward
53         episode_reward += r
54         if isdone:
55             time.sleep(0.1)
56             break
57         episode_rewards.append(episode_reward)
58         epsilon_values.append(agent.epsilon)
59         agent.decay_epsilon()
60
61     if (episode + 1) % 50 == 0:
62         print(
63             "episode:",
64             episode + 1,
65             "episode_reward:",
66             episode_reward,
67             "epsilon:",
68             agent.epsilon,
69         )
70     print('\ntraining over\n')
71
72     # close the render window after training.
73     env.close()
74
75     import matplotlib.pyplot as plt
76     plt.figure(figsize=(12, 6))
77     plt.suptitle('Dyna-Q Algorithm Performance', fontsize=16, fontweight='bold')
78     plt.subplot(1, 2, 1)
79     plt.plot(epsilon_values)
80     plt.title('Epsilon Decay Curve')
81     plt.xlabel('Episode')
82     plt.ylabel('Epsilon')
83
84     def moving_average(data, window_size=20):
85         return np.convolve(data, np.ones(window_size)/window_size, mode='valid')
86
87     avg_rewards = moving_average(episode_rewards)
88     plt.subplot(1, 2, 2)

```

```

89     plt.plot(avg_rewards)
90     plt.title(f'Average_Reward_Curve_(window={20})')
91     plt.xlabel('Episode_(Windowed)')
92     plt.ylabel('Average_Reward')
93
94     plt.tight_layout()
95     plt.show()
96
97
98
99     s = env.reset()
100    agent.epsilon = 0.0
101    episode_reward = 0
102    isdone = False
103    path = [s]
104
105    while not isdone:
106        env.render()
107        time.sleep(0.3)
108        a = agent.choose_action(s)
109        s, r, isdone, info = env.step(a)
110        path.append(s)
111        episode_reward += r
112
113    print(f"Test_complete!_Final_path_reward:{episode_reward}")
114
115    img = env.render(mode='rgb_array')
116    env.close()
117    plt.figure(figsize=(12, 4))
118    plt.imshow(img)
119    coords = [(divmod(s, 12)[1] * img.shape[1] / 12 + img.shape[1] / 24,
120              divmod(s, 12)[0] * img.shape[0] / 4 + img.shape[0] / 8) for s in
121              path]
122    plt.plot([c[0] for c in coords], [c[1] for c in coords], 'r-', linewidth=3,
123            marker='o', markersize=5)
124    plt.title(f'Dyna-Q_Path_(Reward:{episode_reward})', fontsize=14, fontweight='
125            bold')
126    plt.axis('off')
127    plt.savefig('dyna_q_path.png', dpi=150, bbox_inches='tight')
128    plt.show()
129    ##### END CODING HERE #####

```

Listing 3: Source Code for Dyna-Q (cliff_walk_dyna_q.py)

附录 D: Source Code for DQN:

```
1  # -*- coding:utf-8 -*-
2  import argparse
3  import os
4  import random
5  import time
6
7  import gym
8  import numpy as np
9  import torch
10 import torch.nn as nn
11 import torch.nn.functional as F
12 import torch.optim as optim
13 from stable_baselines3.common.buffers import ReplayBuffer
14 from torch.utils.tensorboard import SummaryWriter
15
16 def parse_args():
17     """parse arguments. You can add other arguments if needed."""
18     parser = argparse.ArgumentParser()
19     parser.add_argument("--exp-name", type=str, default=os.path.basename(__file__).rstrip(".py"),
20                         help="the name of this experiment")
21     parser.add_argument("--seed", type=int, default=42,
22                         help="seed of the experiment")
23     parser.add_argument("--total-timesteps", type=int, default=500000,
24                         help="total timesteps of the experiments")
25     parser.add_argument("--learning-rate", type=float, default=2.5e-4,
26                         help="the learning rate of the optimizer")
27     parser.add_argument("--buffer-size", type=int, default=100000,
28                         help="the replay memory buffer size")
29     parser.add_argument("--gamma", type=float, default=0.995,
30                         help="the discount factor gamma")
31     parser.add_argument("--target-network-frequency", type=int, default=500,
32                         help="the timesteps it takes to update the target network")
33     parser.add_argument("--batch-size", type=int, default=256,
34                         help="the batch size of sample from the replay memory")
35     parser.add_argument("--start-e", type=float, default=1.0,
36                         help="the starting epsilon for exploration")
37     parser.add_argument("--end-e", type=float, default=0.01,
38                         help="the ending epsilon for exploration")
39     parser.add_argument("--exploration-fraction", type=float, default=0.15,
40                         help="the fraction of `total-timesteps` it takes from start-e to go end-e")
41     parser.add_argument("--learning-starts", type=int, default=5000,
```

```

42         help="timestep_to_start_learning")
43     parser.add_argument("--train-frequency", type=int, default=4,
44         help="the frequency of training")
45     parser.add_argument("--max-grad-norm", type=float, default=10.0,
46         help="the maximum norm for gradient clipping")
47     args = parser.parse_args()
48     args.env_id = "LunarLander-v2"
49     return args
50
51 def make_env(env_id, seed):
52     """construct the gym environment"""
53     env = gym.make(env_id)
54     env = gym.wrappers.RecordEpisodeStatistics(env)
55     env.seed(seed)
56     env.action_space.seed(seed)
57     env.observation_space.seed(seed)
58     return env
59
60 class QNetwork(nn.Module):
61     """
62     comments:
63
64     the neural network model for approximating Q value function
65
66     Inputs: State
67     Outputs: Q-values for each possible action in that state.
68
69     Here:
70     Input layer: 8 (state dimension)
71     Hidden layer 1: 120 neurons, ReLU activation
72     Hidden layer 2: 84 neurons, ReLU activation
73     Output layer: 4 (action dimension)
74     """
75     def __init__(self, env):
76         super().__init__()
77         self.network = nn.Sequential(
78             nn.Linear(np.array(env.observation_space.shape).prod(), 120),
79             nn.ReLU(),
80             nn.Linear(120, 84),
81             nn.ReLU(),
82             nn.Linear(84, env.action_space.n),
83         )
84
85     def forward(self, x):
86         return self.network(x)

```

```

87
88 def linear_schedule(start_e: float, end_e: float, duration: int, t: int):
89     """
90     comments:
91
92     Implements a linear decay for epsilon ( ) as part of the -greedy strategy.
93     - start_e: The initial value of epsilon
94     - end_e: The final value of epsilon
95     - duration: The total number of timesteps over which to decay from start_e to
96                 end_e.
97     - t: The current timestep.
98
99     When t >= duration, epsilon will be equal to end_e.
100    When t < duration, epsilon will decrease linearly from start_e to end_e.
101    """
102    slope = (end_e - start_e) / duration
103    return max(slope * t + start_e, end_e)
104
105 if __name__ == "__main__":
106
107     """parse the arguments"""
108     args = parse_args()
109     run_name = f"{args.env_id}_{args.exp_name}_{args.seed}_{int(time.time())}"
110
111     """we utilize tensorboard to log the training process"""
112     writer = SummaryWriter(f"runs/{run_name}")
113     writer.add_text(
114         "hyperparameters",
115         "|param|value|\n|-|\n%s" % ("\n".join([f"|{key}|{value}|" for key, value
116                                                in vars(args).items()])),
117     )
118
119     """
120     comments:
121
122     set the random seed to make the experiment reproducible(for numpy, torch, gym,
123         random)
124
125     check whether cuda is available, if available, use cuda to accelerate the
126         training, else use cpu
127     """
128
129     random.seed(args.seed)
130     np.random.seed(args.seed)
131     torch.manual_seed(args.seed)
132     torch.backends.cudnn.deterministic = True
133     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

```

```

128     """
129     comments:
130     create the gym environment
131     envs: the vectorized environment
132     """
133     envs = make_env(args.env_id, args.seed)
134
135     """
136     comments:
137     Initialize the Q-network, target network, and optimizer.
138
139     The target network is a copy of the Q-network and is used to stabilize
140     training.
141     - q_network (Main): Used for action selection (exploitation) and is the
142       network that gets updated by the optimizer.
143     - target_network: Used to calculate the TD Target value. Its weights are
144       frozen. It stabilizes training.
145     - Adam optimizer: Train the parameters of the q_network.
146     """
147     q_network = QNetwork(envs).to(device)
148     optimizer = optim.Adam(q_network.parameters(), lr=args.learning_rate)
149     target_network = QNetwork(envs).to(device)
150     target_network.load_state_dict(q_network.state_dict())
151
152     """
153     comments:
154     Initialize the Replay Buffer.
155     - DQN uses a replay buffer to store past experiences (s, a, r, s', d).
156     Sampling random batches from the buffer breaks the correlation between
157     consecutive samples, stabilizing training.
158     """
159     rb = ReplayBuffer(
160         args.buffer_size,
161         envs.observation_space,
162         envs.action_space,
163         device,
164         handle_timeout_termination=False,
165     )
166
167     """
168     comments:
169     Initialize the environment by resetting it and getting the first observation (
170     state).
171     """
172     obs = envs.reset()

```

```

168     for global_step in range(args.total_timesteps):
169
170         """
171         comments:
172         Calculate the current value of Epsilon.
173         Uses the linear_schedule function based on the current global_step to
174             determine the probability of exploration.
175         """
176         epsilon = linear_schedule(args.start_e, args.end_e, args.
177             exploration_fraction * args.total_timesteps, global_step)
178
179         """
180         comments:
181         Epsilon-Greedy ( -greedy) Action Selection.
182         - With probability 'epsilon', choose a random action (exploration).
183         - With probability '1-epsilon', choose the action with the highest
184             predicted Q-value from the main q_network (exploitation).
185         """
186         if random.random() < epsilon:
187             actions = envs.action_space.sample()
188         else:
189             q_values = q_network(torch.Tensor(obs).to(device))
190             actions = torch.argmax(q_values, dim=0).cpu().numpy()
191
192         """
193         comments:
194         Interact with the environment.
195         Inputs:
196             - actions: The actions chosen by the agent based on the -greedy
197                 strategy.
198         Outputs:
199             - next_obs (s'): The next state.
200             - rewards (r): The immediate reward.
201             - dones (d): whether the episode has ended.
202             - infos: extra info .
203         """
204         next_obs, rewards, dones, infos = envs.step(actions)
205         # envs.render() # close render during training
206
207         if dones:
208             print(f"global_step={global_step}, episodic_return={infos['episode']['r']}")
209             writer.add_scalar("charts/episodic_return", infos["episode"]["r"],
210                 global_step)
211             writer.add_scalar("charts/episodic_length", infos["episode"]["l"],

```

```

        global_step)
207
208     """
209     comments:
210     Store the transition (s, a, r, s', d) in the replay buffer.
211     """
212     rb.add(obs, next_obs, actions, rewards, dones, infos)
213
214     """
215     comments:
216     Update the current observation to the next observation.
217     """
218     obs = next_obs if not dones else envs.reset()
219
220     if global_step > args.learning_starts and global_step % args.
        train_frequency == 0:
221
222         """
223         comments:
224         Sample a random batch of experiences from the Replay Buffer.
225         """
226         data = rb.sample(args.batch_size)
227
228         """
229         comments:
230         Calculate the TD Target value using Double DQN to reduce
            overestimation.
231
232         Standard DQN:  y_j = r +      * max_a' Q_target(s', a')
233         Double DQN:   y_j = r +      * Q_target(s', argmax_a' Q(s', a'))
234
235         Double DQN decouples action selection and evaluation:
236         - Use main Q-network to SELECT the best action
237         - Use target network to EVALUATE that action
238         This reduces the overestimation bias of standard DQN.
239         """
240         with torch.no_grad():
241             # use online network to select actions
242             next_q_values = q_network(data.next_observations)
243             next_actions = next_q_values.argmax(dim=1, keepdim=True)
244             # Use target network to evaluate the selected actions
245             next_q_target_values = target_network(data.next_observations)
246             target_max = next_q_target_values.gather(1, next_actions).squeeze
                ()
247

```



```

248         td_target = data.rewards.flatten() + args.gamma * target_max * (1
                - data.dones.flatten())
249     old_val = q_network(data.observations).gather(1, data.actions).squeeze
        ()
250     loss = F.mse_loss(td_target, old_val)
251
252     """
253     comments:
254     Log the loss and average Q-value to TensorBoard to monitor the
        training process.
255     """
256     if global_step % 100 == 0:
257         writer.add_scalar("losses/td_loss", loss, global_step)
258         writer.add_scalar("losses/q_values", old_val.mean().item(),
                global_step)
259
260     """
261     comments:
262     Perform backpropagation and update the network.
263     Apply gradient clipping to prevent gradient explosion.
264     """
265     optimizer.zero_grad()
266     loss.backward()
267     nn.utils.clip_grad_norm_(q_network.parameters(), args.max_grad_norm)
268     optimizer.step()
269
270     """
271     comments:
272     Periodically update the target network to match the weights of the
        main q_network.
273     """
274     if global_step % args.target_network_frequency == 0:
275         target_network.load_state_dict(q_network.state_dict())
276
277     """close the env and tensorboard logger"""
278     envs.close()
279     writer.close()

```

Listing 4: Source Code for DQN (dqn.py)