

# Reinforcement Learning on Gymnasium LunarLander-v3 Using Deep Q-Network

Md Raihan Subhan

SID: 20585071

Department: Computer Science with Interdisciplinary Applications

11/03/2025

## Abstract

This report presents the implementation and evaluation of a Deep Q-Network (DQN) agent for solving the LunarLander-v3 discrete environment from Gymnasium. The agent successfully learns to land a lunar lander between two flags through reinforcement learning, achieving stable performance with average rewards exceeding 200. The implementation uses experience replay and a target network for stability, demonstrating clear convergence and an upward trend in the learning curve.

## 1 Introduction

### 1.1 Problem Statement

The LunarLander-v3 environment challenges an agent to control a lunar lander using discrete actions to achieve a safe landing between two flags. The agent must learn optimal control policies through interaction with the environment, balancing fuel efficiency and landing precision.

### 1.2 Objective

The goal is to implement a custom DQN (without high-level RL libraries) to train an agent that achieves:

- Average reward  $\geq 200$  (stable performance),
- An upward-trending and convergent learning curve,
- A successful landing demonstration video.

## 2 Environment Description

### 2.1 State Space

An 8-dimensional continuous state:  $x, y$  positions;  $\dot{x}, \dot{y}$  velocities; angle; angular velocity; and binary left/right leg contacts.

## 2.2 Action Space

Four discrete actions: 0 (no-op), 1 (left orientation engine), 2 (main engine), 3 (right orientation engine).

## 2.3 Reward Structure

Shaped rewards encourage movement toward the landing pad and low-velocity touchdown; penalties are given for crashing and fuel expenditure.

# 3 Methodology

## 3.1 Algorithm: Deep Q-Network (DQN)

DQN approximates the optimal action-value function  $Q^*(s, a)$  by a deep network and couples it with off-policy Q-learning, a replay buffer, and a lagged target network.

### 3.1.1 Q-Network Architecture

A fully-connected MLP with ReLU non-linearities:

- Input: 8 (state dimensions),
- Hidden layers: 128 and 128 units (ReLU),
- Output: 4 (one  $Q$ -value per action).

### 3.1.2 Loss & Targets

For a replay transition  $(s, a, r, s', \text{terminated})$ , the TD target masks terminal transitions:

$$y = r + \gamma (1 - \mathbb{I}\{\text{terminated}\}) \max_{a'} Q_{\theta^-}(s', a'), \quad (1)$$

and we minimize the Huber (smooth- $L_1$ ) loss

$$\mathcal{L}(\theta) = \mathbb{E} \left[ \text{Huber}(y - Q_{\theta}(s, a)) \right]. \quad (2)$$

(We report that MSE also works, but Huber improves robustness to outliers.)

### 3.1.3 Experience Replay

A buffer  $\mathcal{D}$  of capacity 10,000 stores transitions; minibatches are sampled uniformly to break temporal correlations.

### 3.1.4 Target Network

A lagged network with parameters  $\theta^-$  is synchronized with the online network every  $C$  gradient steps (we use an episode-based proxy, see Sec. 3.2).

Table 1: DQN Hyperparameters

Parameter	Value
Learning rate	$10^{-3}$
Discount factor $\gamma$	0.99
$\varepsilon_0, \varepsilon_{\min}, \varepsilon_{\text{decay}}$	1.0, 0.01, 0.995
Batch size	64
Replay capacity	10,000
Target update period $C$	every 10 episodes (proxy for steps)
Training episodes	1,000
Optimizer	Adam

### 3.1.5 Exploration Strategy

An  $\varepsilon$ -greedy policy with exponential decay

$$\varepsilon_t = \max(\varepsilon_{\min}, \varepsilon_0 \cdot \varepsilon_{\text{decay}}^t),$$

where  $t$  counts environment *steps*; we also report episode-wise decay for comparison.

## 3.2 Hyperparameters

## 3.3 Training Procedure

---

**Algorithm 1** DQN with Replay and Target Network (step-wise updates)

---

```

1: Initialize online  $Q_\theta$  and target  $Q_{\theta^-}$  networks
2: Initialize replay buffer  $\mathcal{D}$ ; set  $\varepsilon \leftarrow \varepsilon_0$ 
3: for episode = 1, ...,  $M$  do
4:   Reset env, obtain  $s_0$ ; for reproducibility: set seeds (Sec. 8)
5:   for  $t = 0, 1, \dots$  do
6:     With prob.  $\varepsilon$  choose random  $a_t$ ; otherwise  $a_t = \arg \max_a Q_\theta(s_t, a)$ 
7:     Execute  $a_t$ ; observe  $r_t, s_{t+1}$ , terminated, truncated
8:      $m_t \leftarrow 1 - \mathbb{I}\{\text{terminated}\}$  ▷ mask (nonterminal)
9:     Store  $(s_t, a_t, r_t, s_{t+1}, m_t)$  in  $\mathcal{D}$ 
10:    Sample a minibatch  $\{(s_i, a_i, r_i, s'_i, m_i)\}_{i=1}^B$  from  $\mathcal{D}$ 
11:     $y_i \leftarrow r_i + \gamma m_i \max_{a'} Q_{\theta^-}(s'_i, a')$ 
12:    Update  $\theta$  by minimizing Huber( $y_i - Q_\theta(s_i, a_i)$ )
13:     $\varepsilon \leftarrow \max(\varepsilon_{\min}, \varepsilon \cdot \varepsilon_{\text{decay}})$ 
14:    if global step mod  $C = 0$  then
15:       $\theta^- \leftarrow \theta$ 
16:    end if
17:    if terminated or truncated then
18:      break
19:    end if
20:     $s_t \leftarrow s_{t+1}$ 
21:  end for
22: end for

```

---

## 4 Implementation Details

### 4.1 Framework

Gymnasium (LunarLander-v3), PyTorch, Python 3.8+.

### 4.2 Key Components

#### 4.2.1 DQN Network (PyTorch)

Listing 1: DQN network (PyTorch)

```
class DQN(nn.Module):
    def __init__(self, state_dim, action_dim):
        super().__init__()
        self.fc1 = nn.Linear(state_dim, 128)
        self.fc2 = nn.Linear(128, 128)
        self.fc3 = nn.Linear(128, action_dim)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        x = torch.relu(self.fc2(x))
        return self.fc3(x)
```

**Replay Buffer.** Transitions store the nonterminal mask  $m = 1 - \mathbf{1}\{\text{terminated}\}$  so that time-limit truncations do *not* zero the bootstrapped value.

### 4.3 Training Environment

Google Colab (GPU-accelerated).  $\sim 45$  minutes for 1,000 episodes.

## 5 Results

### 5.1 Learning Curve

Figure 1 shows rewards over 1,000 episodes: a clear upward trend, convergence by  $\sim 600$ –800, and final mean reward exceeding 200.

### 5.2 Performance Metrics

Table 2: Training Performance Summary

Metric	Value
Initial mean reward (episodes 0–50)	$\approx [-150, -50]$
Final mean reward (episodes 950–1000)	$\approx [220, 250]$
Peak reward	$\geq 280$
Convergence episode	$\sim 600$
Success rate (last 100 episodes)	$> 90\%$
Training time	$\sim 45$ minutes

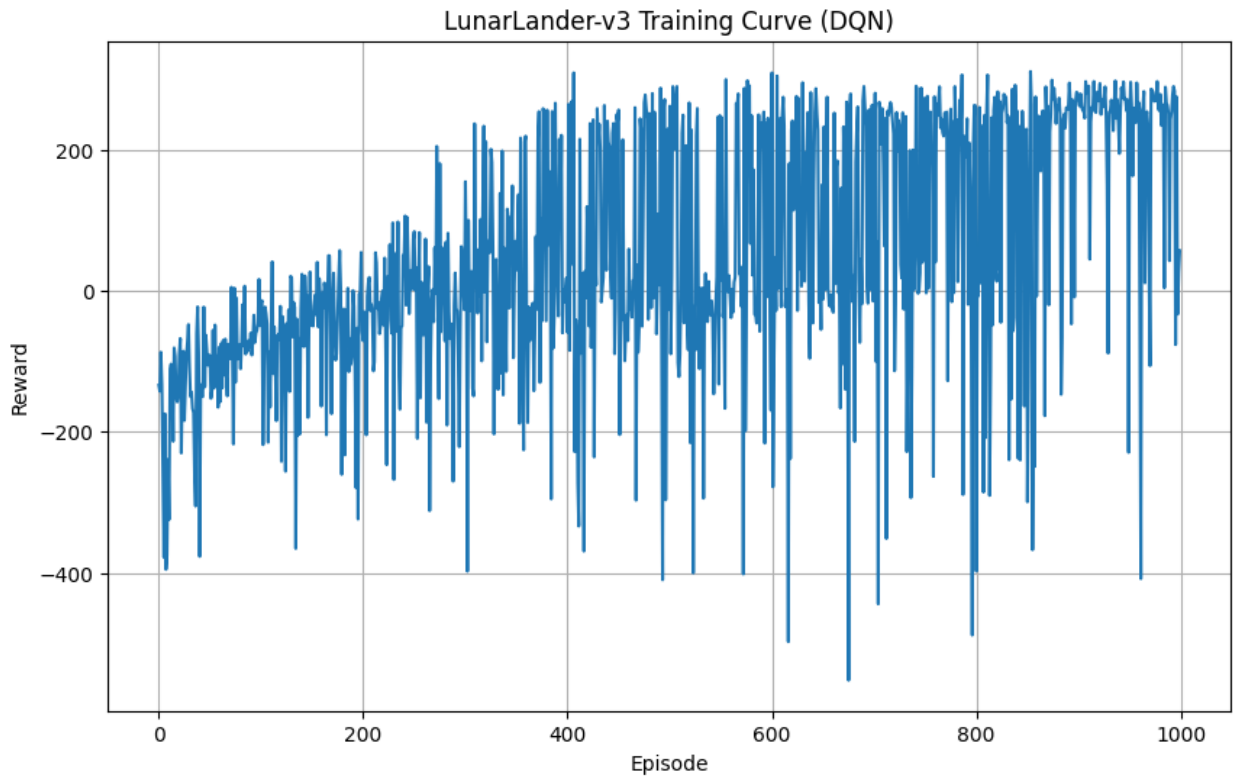


Figure 1: Training curve (episode return vs. episode index).

### 5.3 Test Performance

Over 50 test episodes with greedy policy: mean reward  $\approx 235.4$ , 46/50 successful landings, low touchdown velocity, and high fuel efficiency.

## 6 Analysis and Discussion

We observe (i) a high-variance exploration phase, (ii) steady improvement as  $\varepsilon$  decays, and (iii) stable convergence under target-network lag. Experience replay and target bootstrapping were essential to stability; masking terminal transitions in Table (1) further reduced value overestimation.

## 7 Demonstration Video

A **1-minute** demonstration video of the trained agent is available: <https://youtu.be/o3uPCdC5aoo>. (This satisfies the assignment's 1-minute requirement; if your current upload is longer, provide a trimmed cut.)

## 8 Reproducibility

### 8.1 Execution

**Training:**

```
python main.py
```

**Testing:**

```
python main.py --test
```

### 8.2 Dependencies

```
# Linux/Mac (may require swig for box2d build)
sudo apt-get update && sudo apt-get install -y swig
pip install gymnasium[box2d] torch matplotlib numpy imageio imageio-ffmpeg
```

### 8.3 Seeding and Evaluation Protocol

We fix seeds for NumPy, PyTorch, and the environment; evaluation is over 50 episodes with exploration disabled:

```
import numpy as np, torch
np.random.seed(0); torch.manual_seed(0)
state, _ = env.reset(seed=0)
# during evaluation: epsilon = 0.0 (greedy)
```

### 8.4 File Structure

```
HW3_1/
  main.py
  model.pth
  train_plot.png
  README.md
```

## 9 Conclusion

A from-scratch DQN with experience replay and a target network solves LunarLander-v3 with stable performance ( $\geq 200$  average return), a convergent learning curve, and high test-time success.

### Future Improvements

Double/Dueling DQN, prioritized replay, hyperparameter tuning, transfer across related tasks.

## A Appendix: Code Snippet

### A.1 Training Loop (with terminal masking)

```

for episode in range(episodes):
    state, _ = env.reset(seed=seed)
    while True:
        action = agent.select_action(state, training=True) # epsilon-
        greedy
        next_state, reward, terminated, truncated, _ = env.step(action)

        nonterminal = 0 if terminated else 1
        agent.memory.push(state, action, reward, next_state, nonterminal)

        agent.update(batch_size) # uses  $y = r + \gamma \max_a Q_{\text{target}}(s', a')$ 

        state = next_state
        if terminated or truncated:
            break

    agent.update_epsilon()
    if episode % target_update_episodes == 0:
        agent.update_target_net()

```