

---

# Deep Learning Lab5 Report

陳以瑄 (ID: 313551047)

---

## 1 Introduction

In this lab, I implemented the vanilla DQN and three enhancement techniques: (1) Double DQN, (2) Prioritized Experience Replay buffer, and (3) multi-step return, to solve the CartPole and Pong Atari tasks. Although the vanilla DQN is sufficient to solve the CartPole problem, it requires a significant amount of time to solve the Pong problem, taking about 4 million environment steps. By applying the enhancement techniques, the Pong problem can be solved in just 550k environment steps. I also conducted an ablation study on the CartPole environment, demonstrating that each of the three methods individually improves performance.

## 2 Implementation details

### 2.1 Task1 - Vanilla DQN on CartPole

#### 2.1.1 Neural Network Architecture

For Task 1, aimed at solving the CartPole task, I implemented a simple neural network architecture. The network consists of two hidden layers, each with a dimension of 32.

```
1 class DQN(nn.Module):
2     def __init__(self, input_dim, hidden_dim, num_actions):
3         super(DQN, self).__init__()
4         self.network = nn.Sequential(
5             nn.Linear(input_dim, hidden_dim),
6             nn.ReLU(),
7             nn.Linear(hidden_dim, hidden_dim),
8             nn.ReLU(),
9             nn.Linear(hidden_dim, num_actions)
10        )
11    def forward(self, x):
12        return self.network(x)
```

#### 2.1.2 Replay Buffer

Since this is a vanilla DQN, the replay buffer uses uniform random sampling. I store transitions in a deque and sample batches randomly.

```

1 self.memory = deque(maxlen=args.memory_size)
2
3 batch = random.sample(self.memory, self.batch_size)
4 states, actions, rewards, next_states, dones = zip(*batch)

```

### 2.1.3 Bellman Error for DQN

The DQN loss is defined as  $\mathcal{L}_{\text{DQN}}(\theta) = \sum_{(s,a,r,s') \in D} (r + \gamma \max_{a' \in A} Q(s', a'; \bar{\theta}) - Q(s, a; \theta))^2$  where  $\theta$  are the parameters of the current Q-network, and  $\bar{\theta}$  are the parameters of the target network.

**Line 2:** Predict the current Q-value  $Q(s, a; \theta)$  using the Q-network.

**Line 4:** Obtain the maximum next Q-value  $\max_{a' \in A} Q(s', a'; \bar{\theta})$  using the target network.

**Line 5:** Compute the target Q-value as  $r + \gamma \max_{a'} Q(s', a'; \bar{\theta})$ .

**Lines 6–7:** Calculate the Bellman error (loss) as the Mean Squared Error (MSE) between the predicted Q-values and the target Q-values.

```

1 # After getting states, next_states, actions, rewards, dones
2 q_values = self.q_net(states).gather(1, actions.unsqueeze(1)).squeeze(1)
3 with torch.no_grad():
4     q_nexts = self.target_net(next_states).max(1)[0]
5     q_targets = rewards + self.gamma * q_nexts * (1 - dones)
6 criterion = nn.MSELoss()
7 loss = criterion(q_values, q_targets)
8
9 self.optimizer.zero_grad()
10 loss.backward()
11 self.optimizer.step()

```

## 2.2 Task2 - Vanilla DQN on Atari Pong

### 2.2.1 Neural Network Architecture

Since Task 2 involves solving a more complex environment, I design a deeper neural network consisting of two parts: a convolutional feature extractor with three layers, and a fully connected classifier that maps the extracted features to action values.

```

1 class DQN(nn.Module):
2     def __init__(self, num_actions):
3         super(DQN, self).__init__()
4         self.cnn = nn.Sequential(nn.Conv2d(4, 32, kernel_size=8, stride=4),
5                                   nn.ReLU(True),
6                                   nn.Conv2d(32, 64, kernel_size=4, stride=2),
7                                   nn.ReLU(True),
8                                   nn.Conv2d(64, 64, kernel_size=3, stride=1),
9                                   nn.ReLU(True)
10                                )

```

```

11         self.classifier = nn.Sequential(nn.Linear(7*7*64, 512),
12                                         nn.ReLU(True),
13                                         nn.Linear(512, num_actions)
14                                         )
15
16     def forward(self, x):
17         x = x.float() / 255.
18         x = self.cnn(x)
19         x = torch.flatten(x, start_dim=1)
20         x = self.classifier(x)
21         return x

```

## 2.3 Task3 - Enhance DQN on Atari Pong

### 2.3.1 Modify DQN to Double DQN

The Double DQN loss is defined as:

$\mathcal{L}_{\text{DDQN}}(\theta) = \sum_{(s,a,r,s') \in D} (r + \gamma Q(s', \arg \max_{a' \in A} Q(s, a'; \bar{\theta}) - Q(s, a; \theta))^2$ , where  $\theta$  are the parameters of the current Q-network, and  $\bar{\theta}$  are the parameters of the target network.

**Line 2:** Predict the current Q-value  $Q(s, a; \theta)$  using the behavior net.

**Line 4:** Use the behavior network to select the next action  $\max_{a'} Q(s', a'; \theta)$ .

**Line 5:** Evaluate the next Q-value  $Q(s', \arg \max_{a'} Q(s, a'; \bar{\theta}))$  using the target network to avoid overestimation bias.

**Line 6—7:** Compute the Bellman error (loss) as the Mean Squared Error (MSE) between the predicted Q-values and the target Q-values.

```

1 # After getting states, next_states, actions, rewards, dones
2 q_values = self.q_net(states).gather(1, actions.unsqueeze(1)).squeeze(1)
3 with torch.no_grad():
4     next_actions = torch.argmax(self.q_net(next_states), dim = 1)
5     q_nexts = self.target_net(next_states.squeeze(-1)).gather(1,
6     next_actions.unsqueeze(1)).squeeze(1)
7     q_targets = rewards+ self.gamma * q_nexts*(1-dones)
8
9 criterion = nn.MSELoss()
10 loss = criterion(q_values, q_targets)
11
12 self.optimizer.zero_grad()
13 loss.backward()
14 self.optimizer.step()

```

### 2.3.2 Implement the memory buffer for PER

1. Add

Use a First-In-First-Out (FIFO) approach to add transitions.

Assign priorities  $p$  as  $p = |\delta| + \epsilon$ , and save  $p^\alpha$  for efficient future computation.

```

1 def add(self, transition, error=1):
2     if len(self.buffer) < self.capacity:
3         self.buffer.append(None)
4     self.buffer[self.pos] = transition
5     self.priorities[self.pos] = (abs(error) + self.epsilon)** self.alpha
6     self.pos = (self.pos+1)% self.capacity

```

In DQNAgent run(), after getting the (s,a,r,s',a'), insert the transition into the PER buffer with the highest priority.

```

1 if len(self.memory.buffer):
2     prior = self.memory.priorities[: len(self.memory.buffer)].max()
3     self.memory.add((state, action, reward, next_state, done), error=prior)
4 else:
5     self.memory.add((state, action, reward, next_state, done), error = 1)

```

## 2. Sample

Sample a batch of transitions with the probability  $P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$ . Then, calculate the importance sampling weight  $w_i = \left(\frac{1}{NP(i)}\right)^\beta$  and return the indices, samples, and weights.

```

1 def sample(self, batch_size):
2     if len(self.buffer) == self.capacity:
3         priorities = self.priorities
4     else:
5         priorities = self.priorities[: self.pos]
6     probs = priorities/priorities.sum()
7     indices = np.random.choice(len(probs), batch_size, p = probs)
8     samples = [self.buffer[idx] for idx in indices]
9
10    weights = (len(self.buffer)* probs[indices])**(-self.beta)
11    weights/= weights.max()
12    return indices, samples, weights

```

In DQNAgent train(), sample a batch of transition.

```

1 indices, batch, weights = self.memory.sample(self.batch_size)
2 states, actions, rewards, next_states, dones = zip(*batch)

```

## 3. Update

```

1 def update_priorities(self, indices, errors):
2     for idx, err in zip(indices, errors):
3         self.priorities[idx] = (abs(err) + 1e-6) ** self.alpha
4     return

```

In DQNAgent train(), after getting the new Bellman error, update the priorities.

```

1 BM_errors = (q_values - q_targets.squeeze()).detach().cpu().numpy()
2 self.memory.update_priorities(indices, BM_errors)

```

### 2.3.3 Modify the 1-step return to multi-step return

For each step, save the tuple  $(s, a, r, s', d)$  into the `nstep_buffer`. Once the buffer reaches `n`-step, calculate the discounted cumulative `n`-step reward and save  $(s, a, n\text{-step } R, n\text{-step } s', n\text{-step } d)$  into the replay buffer.

```
1 self.nstep_buffer.append( (state, action, reward, next_state, done))
2 if len(self.nstep_buffer) == self.n_steps:
3     R = 0.0
4     s_n = next_state
5     done_n = done
6     for k, (_, _, r, s_next, d) in enumerate(reversed(self.nstep_buffer)):
7         R = r + (self.gamma * R * (1 - d))
8         if d:
9             s_n = s_next
10            done_n = True
11            break
12    s_0, a_0, _, _, _ = self.nstep_buffer[0]
13    self.memory.append((s_0, a_0, R, s_n, done_n))
```

Since it is a `N`-step target, when we compute the `q_target`, we should discount by  $\gamma^N$

```
1 q_targets = rewards+ (self.gamma ** self.n_steps) * q_nexts*(1-dones)
```

## 2.4 Explain how you use Weight & Bias to track the model performance

I just followed the code given in the Sample Code and added the log in the `DQNAgent`'s `train()` to log the training loss.

```
1 wandb.log({
2     "Train Loss": loss.item(),
3     "Q Mean": q_values.mean().item(),
4     "Q Std": q_values.std().item()
5 })
```

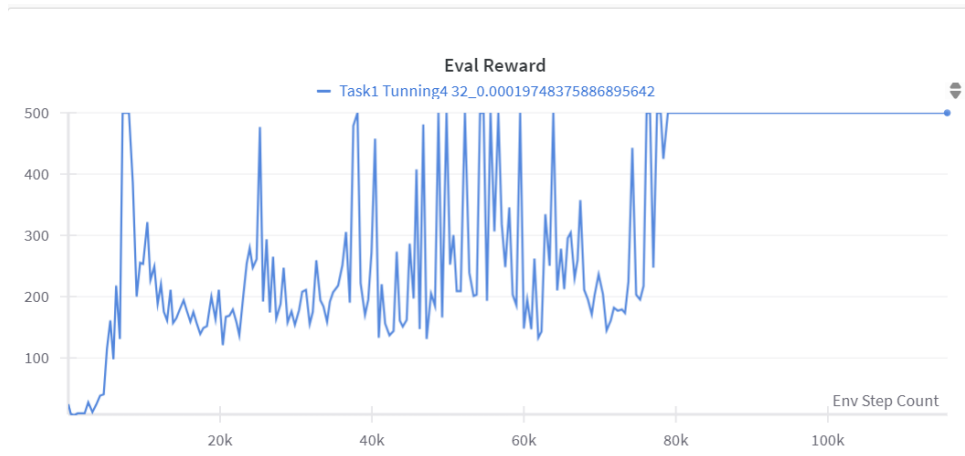
## 3 Analysis & Discussion

### 3.1 Training Curve

#### 3.1.1 Task 1

Hyperparameter: `lr = 0.00019`; `batch size = 128`; `replay-start-size = 500`

Result: Converge in 80k env steps.



### 3.1.2 Task 2

hyperparameter: lr = 0.0001; batch size = 32; replay-start-size = 50000

Result: Reach 500 eval reward in 4M env step



### 3.1.3 Task 3

hyperparameter: n-step = 4; lr = 0.00025; batch size = 32; replay-start-size = 30000

Result: Reach 500 eval reward around 550k env step



## 3.2 Comparison

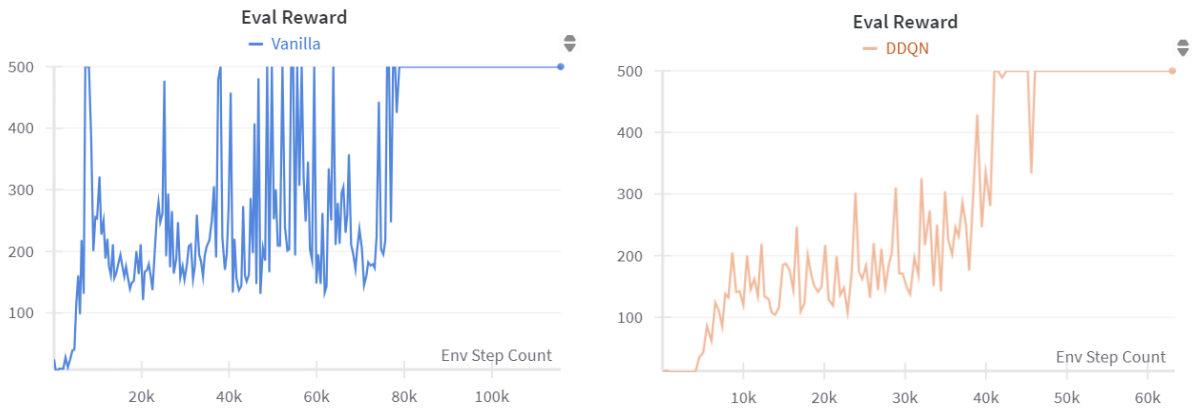
### 3.2.1 Compare the Result over Task 2 and 3

From the plot, we can easily observe that applying the three enhancement techniques significantly improves overall performance. The original version required 4 million environment steps to reach the goal, whereas the enhanced version only took 550k environment steps.



### 3.2.2 Compare the Result between DQN and DDQN on CartPole

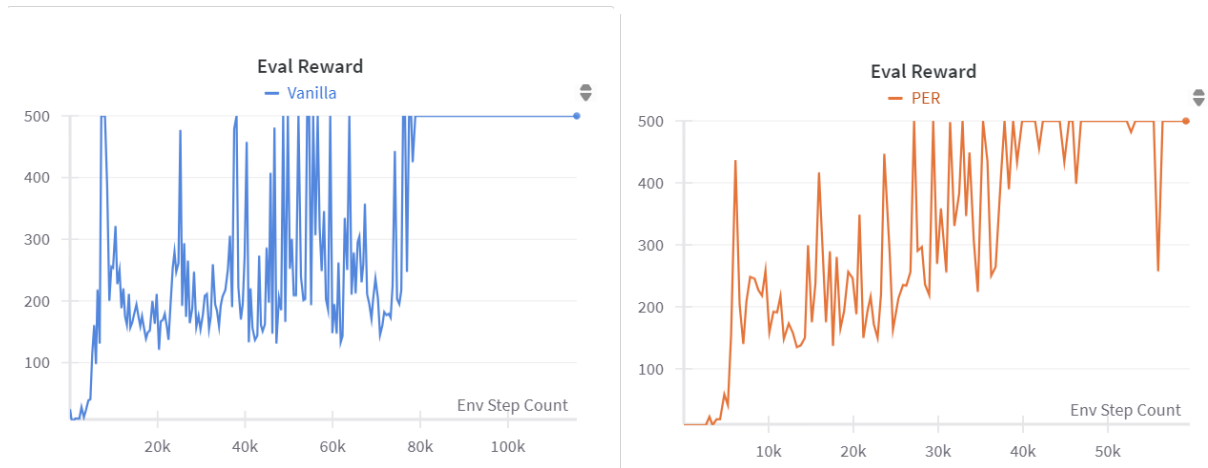
We can observe that the DDQN version exhibits more stable reward values, as it avoids overestimating the Q-values. As a result, the DDQN version converges faster (50k < 80k environment steps).



### 3.2.3 Compare the Result of PER on CartPole

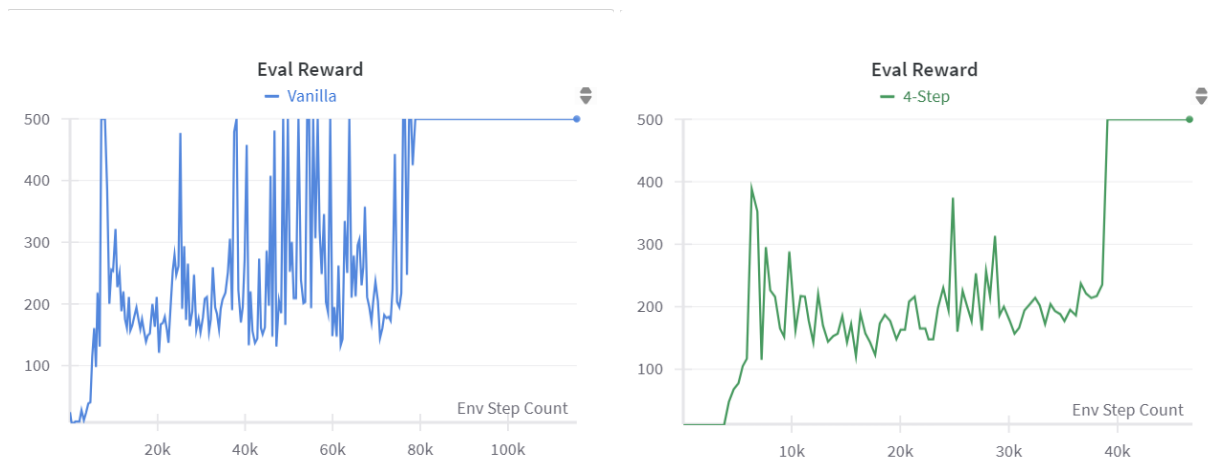
We can observe that, since the PER buffer prioritizes transitions with larger Bellman errors—indicating they require more learning—the version with PER converges faster

(40k < 80k environment steps).



### 3.2.4 Compare the Result of NStep Return on CartPole

We can observe that the version using 4-step returns exhibits more stable reward values, resulting in a steadier increase and faster convergence (50k < 80k environment steps).



## 4 Other training strategy

### 4.1 RMSProp Optimizer

Since the environment is noisy, I applied the RMSProp optimizer instead of Adam to achieve more stable convergence.

### 4.2 Reward Clipping

Since the reward in Pong ranges from -21 to 21, which introduces significant variability, I clip the reward to between  $\pm 1$  to indicate to the agent whether it wins or loses, making the learning process more stable.



### 4.3 Optuna Hyperparameter Tunning

Since the learning rate is critical, a small learning rate can lead to slower convergence, while a large learning rate may cause overshooting. Therefore, I use Optuna to search for the optimal learning rate within the range of  $[1e-5, 1e-3]$ .

