

Deep Q Network (DQN) Implementation

with Atari Breakout

I. Abstract

1.1 Overview

This project implements a Deep Q-Network (DQN) for playing Atari games, specifically focusing on the game "Breakout." The implementation includes enhancements to the basic DQN structure, such as Dueling Network Architectures and Double Q-Learning, which make the agent's learning process more efficient and robust. The experiment setup is designed to run game environments on multiple processes for efficient sampling. Each component of the DQN in this project is well modularized to facilitate understanding and application of DQN in practical reinforcement learning scenarios.

1.2 Model Architecture

The DQN model employs a convolutional neural network (CNN) that processes the state inputs from the environment. These states are typically frames or images that describe the current situation in the environment. The network then outputs Q-values for each possible action, which represent the expected future rewards that can be obtained by taking each action in the current state.

1.3 Dueling Network Architecture

A key feature of the model is the dueling network architecture, which consists of two separate estimators:

1. Value Function $V(s)$: Estimates how good it is to be in a given state s
2. Advantage Function $A(s, a)$: Estimates the advantage of taking a specific action a over others in state s

The final Q-value for each action is computed using the equation:

$$Q(s, a) = V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right)$$

Where $|\mathcal{A}|$ is the number of possible actions. This equation adjusts the advantage values so that they are zero-centered, which stabilizes the training process by reducing the variance of the advantage estimates.

1.4 Training Process

1. **Sampling:** Actions are sampled using an ϵ -greedy policy, where ϵ is gradually reduced. This means the model mostly exploits the best-known actions while occasionally exploring other actions to discover potentially better strategies.
2. **Experience Replay:** Interactions with the environment (state transitions) are stored in a replay buffer. Training samples are drawn from this buffer, which helps to break the correlation between consecutive training samples and stabilizes learning.
3. **Double Q-Learning:** To mitigate the overestimation of Q-values, two networks are maintained: the primary network and a target network. The target network's weights are periodically updated with the weights from the primary network, providing stable targets for training.
4. **Loss Function:** The loss is computed as the difference between the current predicted Q-values and the target Q-values, which are adjusted by the reward received and the discounted highest Q-value of the next state, as predicted by the target network.
5. **Prioritized Experience Replay:** The replay buffer is enhanced with a prioritization mechanism that more frequently samples transitions with high temporal difference (TD) errors, which are indicative of significant learning potential.

By continuously interacting with the environment, updating the network weights, and adjusting the exploration/exploitation balance, the model progressively learns to maximize rewards over time, ideally converging to an optimal policy that dictates the best actions to take from any given state.

II. Implementation Code Core Section Annotation

In this section, I will provide screenshots of the code, along with comments on the core sections of the project. At the beginning of each function, there is a short description (orange color).

2.1 Model.py

2.1.1 Import Section

```
12 import torch
13 from torch import nn
14 from labml_helpers.module import Module
```

2.1.2 Class Model

```
model.py > Model > __init__
18 class Model(Module):
49     def __init__(self):
50         super().__init__()
51         self.conv = nn.Sequential(
52             # First convolutional layer: takes a 84x84 frame with 4 channels as input
53             # and outputs a 20x20 frame with 32 channels
54             nn.Conv2d(in_channels=4, out_channels=32, kernel_size=8, stride=4),
55             nn.ReLU(),
56
57             # Second convolutional layer: reduces the dimension from 20x20 to 9x9 with 64 output channels
58             nn.Conv2d(in_channels=32, out_channels=64, kernel_size=4, stride=2),
59             nn.ReLU(),
60
61             # Third convolutional layer: processes the 9x9 frame to output a 7x7 frame with 64 channels
62             nn.Conv2d(in_channels=64, out_channels=64, kernel_size=3, stride=1),
63             nn.ReLU(),
64         )
65
66         # A fully connected layer takes the flattened frame from third convolution layer,
67         # and outputs 512 features
68         self.lin = nn.Linear(in_features=7 * 7 * 64, out_features=512)
69         self.activation = nn.ReLU()
70
71         # This head gives the state value $V$
72         self.state_value = nn.Sequential(
73             nn.Linear(in_features=512, out_features=256),
74             nn.ReLU(),
75             nn.Linear(in_features=256, out_features=1),
76         )
77         # This head gives the action value $A$
78         self.action_value = nn.Sequential(
79             nn.Linear(in_features=512, out_features=256),
80             nn.ReLU(),
81             nn.Linear(in_features=256, out_features=4),
82         )
```

```
84     def forward(self, obs: torch.Tensor):
85         # Convolution
86         h = self.conv(obs)
87
88         # Reshape for linear layers
89         h = h.reshape((-1, 7 * 7 * 64))
90
91         # Pass the reshaped tensor 'h' through a linear layer followed by an activation function.
92         h = self.activation(self.lin(h))
93
94         # Compute the advantage of each action using the 'action_value' network.
95         action_value = self.action_value(h)
96
97         # Compute the value of being in the given state using the 'state_value' network.
98         state_value = self.state_value(h)
99
100        # Center the action values by subtracting the mean action value.
101        action_score_centered = action_value - action_value.mean(dim=-1, keepdim=True)
102
103        # Calculate the Q-value for each action by combining the state value and the centered action values.
104        q = state_value + action_score_centered
105
106        return q
```

2.2 Experiment.py

2.2.1 Import Section

```
15 import numpy as np
16 import torch
17 from labml import tracker, experiment, logger, monit
18 from labml.internal.configs.dynamic_hyperparam import FloatDynamicHyperParam
19 from labml.helpers.schedule import Piecewise
20 from labml.nn.rl.dqn import QFuncLoss
21 from labml.nn.rl.dqn.model import Model
22 from labml.nn.rl.dqn.replay_buffer import ReplayBuffer
23 from labml.nn.rl.game import Worker
24
25 # Select the computing device
26 if torch.cuda.is_available():
27     device = torch.device("cuda:0")
28 else:
29     device = torch.device("cpu")
30
31
32 def obs_to_torch(obs: np.ndarray) -> torch.Tensor:
33     """Scale observations from `[0, 255]` to `[0, 1]`"""
34     return torch.tensor(obs, dtype=torch.float32, device=device) / 255.
```

2.2.2 Class Trainer

2.2.2.1 Init Function

```
37 class Trainer:
38     def __init__(self, *, updates: int, epochs: int, n_workers: int, worker_steps: int, mini_batch_size: int,
39                 update_target_model: int, learning_rate: FloatDynamicHyperParam):
40
41         self.n_workers = n_workers # Number of parallel worker processes
42         self.worker_steps = worker_steps # Number of steps each worker takes before updating the model
43         self.train_epochs = epochs # Number of epochs to train with each batch of data
44         self.updates = updates # Total number of model updates to perform
45         self.mini_batch_size = mini_batch_size # Size of each training batch
46         self.update_target_model = update_target_model # Interval at which to update the target model
47         self.learning_rate = learning_rate # Dynamic adjustment of learning rate
48
49
50         # exploration as a function of updates
51         self.exploration_coefficient = Piecewise(
52             [
53                 (0, 1.0),
54                 (25_000, 0.1),
55                 (self.updates / 2, 0.01)
56             ], outside_value=0.01)
57
58         # beta for replay buffer as a function of updates
59         self.prioritized_replay_beta = Piecewise(
60             [
61                 (0, 0.4),
62                 (self.updates, 1)
63             ], outside_value=1)
64
65         # Initialize replay buffer with specific capacity (power of 2) and alpha (0.6)
66         self.replay_buffer = ReplayBuffer(2 ** 14, 0.6)
67
68         # Model for sampling and training
69         self.model = Model().to(device)
70
71         # Target model used in Double Q-learning
72         self.target_model = Model().to(device)
```

```
73
74         # Initialize workers for parallel environment interaction
75         self.workers = [Worker(47 + i) for i in range(self.n_workers)]
76
77         # Initial observations from the environment
78         self.obs = np.zeros((self.n_workers, 4, 84, 84), dtype=np.uint8)
79
80         # Initialize the environment in each worker
81         for worker in self.workers:
82             worker.child.send(("reset", None))
83
84         # Get the initial observations from each worker
85         for i, worker in enumerate(self.workers):
86             self.obs[i] = worker.child.recv()
87
88         # Initialize the loss function with a discount factor
89         self.loss_func = QFuncLoss(0.99)
90
91         # Initialize the optimizer with a learning rate
92         self.optimizer = torch.optim.Adam(self.model.parameters(), lr=2.5e-4)
```

2.2.2.2 _Sample_action Function

```
94 def _sample_action(self, q_value: torch.Tensor, exploration_coefficient: float):
95     """
96     Sample an action using an epsilon-greedy strategy where epsilon is the exploration coefficient.
97     """
98
99     # Sampling doesn't need gradients
100     # Disable gradient calculations for efficiency
101     with torch.no_grad():
102         # Sample the action with highest Q-value. This is the greedy action.
103         greedy_action = torch.argmax(q_value, dim=-1)
104         # Uniformly sample and action
105         random_action = torch.randint(0, q_value.shape[-1], (1,))
106         # Whether to chose greedy action (variable) greedy_action: Tensor
107         is_choose_rand = torch.rand(1, device=q_value.device) < exploration_coefficient
108         # Pick the action based on epsilon-greedy strategy `is_choose_rand`
109         return torch.where(is_choose_rand, random_action, greedy_action).cpu().numpy()
```

2.2.2.3 Sample Function

```
111 def sample(self, exploration_coefficient: float):
112     """
113     Collect data samples from the environment using the current policy and update the replay buffer.
114     """
115     # Disable gradient calculations for efficiency
116     with torch.no_grad():
117         # Perform sampling for defined number of worker steps
118         for t in range(self.worker_steps):
119             # Get Q values for the current observation
120             q_value = self.model(obs_to_torch(self.obs))
121             # Sample actions
122             actions = self._sample_action(q_value, exploration_coefficient)
123
124             # Run sampled actions in each worker and update observations
125             for w, worker in enumerate(self.workers):
126                 worker.child.send(("step", actions[w]))
127
128             # Collect information from each worker
129             for w, worker in enumerate(self.workers):
130                 # Get results after executing the actions
131                 next_obs, reward, done, info = worker.child.recv()
132
133                 # Add transition to replay buffer
134                 self.replay_buffer.add(self.obs[w], actions[w], reward, next_obs, done)
135
136                 # update episode information.
137                 # collect episode info, which is available if an episode finished;
138                 # this includes total reward and length of the episode -
139                 # look at `Game` to see how it works.
140                 if info:
141                     tracker.add('reward', info['reward'])
142                     tracker.add('length', info['length'])
143
144             # update current observation
145             self.obs[w] = next_obs
```

2.2.2.4 Train Function

```
147 def train(self, beta: float):
148     """
149     Train the model using samples from the replay buffer and update the model and target networks.
150     """
151     # Perform training for defined number of epochs
152     for _ in range(self.train_epochs):
153         # Sample from priority replay buffer
154         samples = self.replay_buffer.sample(self.mini_batch_size, beta)
155         # Get the predicted Q-value
156         q_value = self.model(obs_to_torch(samples['obs']))
157
158         # Get the Q-values of the next state for [Double Q-learning](index.html).
159         # Gradients shouldn't propagate for these
160         with torch.no_grad():
161             # Double Q-learning Q(s;theta_i)
162             double_q_value = self.model(obs_to_torch(samples['next_obs']))
163             # Target Q-value for Double Q-learning Q(s;theta_i^-)
164             target_q_value = self.target_model(obs_to_torch(samples['next_obs']))
165         # Calculate Temporal Difference (TD) errors (delta) and loss L(theta)
166         td_errors, loss = self.loss_func(q_value,
167                                         q_value.new_tensor(samples['action']),
168                                         double_q_value, target_q_value,
169                                         q_value.new_tensor(samples['done']),
170                                         q_value.new_tensor(samples['reward']),
171                                         q_value.new_tensor(samples['weights']))
172
173         # Calculate priorities for replay buffer, p_i = |delta_i| + epsilon
174         new_priorities = np.abs(td_errors.cpu().numpy()) + 1e-6
175         # Update replay buffer priorities
176         self.replay_buffer.update_priorities(samples['indexes'], new_priorities)
177
178         # Set learning rate
179         for pg in self.optimizer.param_groups:
180             pg['lr'] = self.learning_rate()
181
182         # Zero out the previously calculated gradients
183         self.optimizer.zero_grad()
184         # Calculate gradients
185         loss.backward()
186         # Clip gradients
187         torch.nn.utils.clip_grad_norm_(self.model.parameters(), max_norm=0.5)
188         # Update parameters based on gradients
189         self.optimizer.step()
```


2.2.2.5 Run_training_loop Function

```
190 def run_training_loop(self):
191     """
192     Execute the main training loop, managing exploration, training, and updates to the target network.
193     """
194
195     # Last 100 episode information
196     tracker.set_queue('reward', 100, True)
197     tracker.set_queue('length', 100, True)
198
199     # Copy to target network initially
200     self.target_model.load_state_dict(self.model.state_dict())
201     # Main update loop
202     for update in monit.loop(self.updates):
203         # Calculate current exploration rate epsilon
204         exploration = self.exploration_coefficient(update)
205         tracker.add('exploration', exploration)
206         # Calculate current beta for the replay buffer beta
207         beta = self.prioritized_replay_beta(update)
208         tracker.add('beta', beta)
209
210         # Sample with current policy
211         self.sample(exploration)
212
213         # Start training after the buffer is full
214         if self.replay_buffer.is_full():
215             # Train the model
216             self.train(beta)
217
218             # Update the target network periodically
219             if update % self.update_target_model == 0:
220                 self.target_model.load_state_dict(self.model.state_dict())
221
222     # Save tracked indicators.
223     tracker.save()
224     # Add a new line to the screen periodically
225     if (update + 1) % 1_000 == 0:
226         logger.log()
```

2.2.2.6 Destroy Function

```
228 def destroy(self):
229     """
230     Clean up resources and stop worker processes at the end of training.
231     """
232     # Send close signal to each worker
233     for worker in self.workers:
234         worker.child.send(("close", None))
```

2.2.3 Main Function

```
237 def main():
238     # Create the experiment
239     experiment.create(name='dqn')
240
241     # Configurations
242     configs = {
243         # Number of updates
244         'updates': 1_000_000,
245         # Number of epochs to train the model with sampled data.
246         'epochs': 8,
247         # Number of worker processes
248         'n_workers': 8,
249         # Number of steps to run on each process for a single update
250         'worker_steps': 4,
251         # Mini batch size
252         'mini_batch_size': 32,
253         # Target model updating interval
254         'update_target_model': 250,
255         # Learning rate.
256         'learning_rate': FloatDynamicHyperParam(1e-4, (0, 1e-3)),
257     }
258
259     # Apply configurations to the experiment
260     experiment.configs(configs)
261
262     # Initialize the trainer
263     m = Trainer(**configs)
264     # Run and monitor the experiment
265     with experiment.start():
266         m.run_training_loop()
267     # Stop the workers
268     m.destroy()
269
270
271 # Run the main function
272 if __name__ == "__main__":
273     main()
```

Reference

Labmlai. (n.d.). *annotated_deep_learning_paper_implementations/labml_nn/rl/dqn* at *master* · *labmlai/annotated_deep_learning_paper_implementations*. GitHub.
https://github.com/labmlai/annotated_deep_learning_paper_implementations/tree/master/labml_nn/rl/dqn