

Deep Q-Network Algorithm

Asfahan Shah (aks7824@psu.edu)

1 Abstract

A simple (Deep Q-Network) DQN is implemented. This algorithm is tested on CartPole environment. The DQN algorithm combines deep neural networks with Q-learning to learn optimal policies for an environment. The implementation has following key components:

- Experience Replay Memory: Stores state transitions. During training we get a randomly sampled batch from experience replay memory. The benefit is that training in a batch is more effective rather than training one by one sample. Another benefit is by randomly sampling we are making data kind of I.I.D (Identically Independently Distributed) which is beneficial for us .
- Target Network: Uses a separate network for generating target Q-values. This network usually remains fixed and is updated via soft updates to provide stable training targets.
- Epsilon-Greedy Exploration: When choosing a action, if we always choose an action based on maximum Q-value it could be detrimental. This is because model is not know here so we have to explore. To do this we employ epsilon-greedy approach. For some probability we choose random action while for other we choose maximum Q-value action. This helps in balancing exploration and exploitation.

The training process follows these steps:

1. Initialize policy and target networks with identical weights
2. For each episode:
 - Reset environment and get initial state
 - For each time step:
 - Select action using epsilon-greedy approach
 - Execute action and observe reward and next state
 - Store transition in replay memory
 - Sample random batch from replay memory
 - Perform Q-learning update on policy network
 - Perform soft update on target network

2 Implementation Analysis

2.1 Policy Architecture

```
1 class Policy(torch.nn.Module):
2     def __init__(self, state, action):
3         super().__init__()
4         # Three-layer neural network
5         self.input = Linear(state, 128)      # First layer maps state space to
6         hidden layer
7         self.dense = Linear(128, 128)        # Hidden layer for feature extraction
```

```

7         self.output = Linear(128, action)    # Output layer produces Q-values for
            each action
8
9     def forward(self, tensor):
10         x = F.relu(self.input(tensor))        # Apply ReLU to first layer output
11         x = F.relu(self.dense(x))            # Apply ReLU to hidden layer output
12         return self.output(x)                # Return raw Q-values

```

The neural network architecture implements a three-layer feed-forward network for Q-value approximation:

- Input Layer: Maps the state space to 128 hidden units.
- Hidden Layer: Maintains representation with 128 units for feature extraction.
- Output Layer: Produces Q-values for each possible action.

ReLU activation functions are used between layers.

2.2 Experience Replay Buffer

```

1 class memory():
2     def __init__(self, capacity):
3         # Initialize circular buffer with fixed capacity
4         self.mem = deque([], maxlen=capacity)
5
6     def add(self, state, action, reward, next_state):
7         # Store transition tuple in memory
8         self.mem.append((state, action, reward, next_state))
9
10    def get_sample(self, size):
11        # Randomly sample batch of a particular size of transitions for training
12        return random.sample(self.mem, size)
13    def __len__(self):
14        return len(self.mem)

```

The experience replay buffer is crucial for stable learning. The following are key points:

- Uses a circular buffer (deque) with fixed capacity.
- Stores transitions as tuples: (state, action, reward, next_state).
- Random sampling breaks temporal correlations in the training data.
- Fixed capacity prevents memory issues.

2.3 Target Network Update Mechanism

```

1 def soft_update(self):
2     # Get current weights of both networks
3     w_1 = self.Policy_net.state_dict()
4     w_2 = self.Target_net.state_dict()
5     # Perform soft update.
6     for key in w_1:
7         # Update formula: theta_target = theta*policy + (1-theta)*target
8         w_2[key] = self.theta * w_1[key] + (1-self.theta) * w_2[key]
9     # Load updated weights into target network
10    self.Target_net.load_state_dict(w_2)

```

2.3.1 Traditional Hard Updates vs Soft Updates

Problems with Hard Updates: In traditional DQN implementations, the target network is updated by completely copying the policy network's weights after every K steps. This approach has following problems:

- **Abrupt Changes:** When the target network is updated, there's a sudden large change in the target Q-values
- **Training Instability:** These abrupt changes can lead to significant oscillations in the learning process
- **Frequency Dilemma:**
 - If K is too small: Target values change too rapidly, leading to instability
 - If K is too large: Target network becomes outdated, slowing learning

Advantages of Soft Updates: Our implementation uses soft updates after each step with a small update rate ($\theta = 0.005$):

- **Gradual Updates:** Instead of copying all weights, we blend a small portion of the policy network's weights into the target network
- **Update Formula:** For each weight w: $w_{target} = \theta \cdot w_{policy} + (1 - \theta) \cdot w_{target}$ where $\theta = 0.005$
- **Continuous Learning:** Updates can occur after every step without destabilizing training
- **Smooth Transitions:** Prevents sudden jumps in target Q-values

2.4 Action Selection Strategy

```
1 def choose_action(self, state):
2     self.steps += 1
3     # Calculate exploration probability using exponential decay
4     threshold = self.END + (self.START - self.END) * \
5         math.exp(-1. * self.steps / self.DECAY)
6     sample = random.random()
7
8     if sample > threshold:
9         # Exploit: choose action with highest Q-value
10        with torch.no_grad():
11            return self.Policy_net(state).max(1)[1].view(1, 1)
12    else:
13        # Explore: choose random action
14        return torch.tensor([[self.env.action_space.sample()]],
15                             device=self.DEVICE, dtype=torch.long)
```

The action selection implements epsilon-greedy exploration. Following are the key points:

- Exploration probability decays exponentially with steps.
- The exploration probability is decaying with steps because after lets say after some large steps, we would want policy to exploit as ideally we would have explored sufficiently.
- During exploration: Random action sampling from action space.
- During exploitation: Selects action with highest Q-value.

2.5 Training Process

```
1 def train(self):
2     # Skip if not enough samples in buffer
3     if len(self.buffer) < self.BATCH_SIZE:
4         return None
5
6     # Sample random batch from replay buffer
7     batch = self.buffer.get_sample(self.BATCH_SIZE)
8     # Create mask for non-final states
9     non_final_index = torch.tensor(
10         tuple([True if i[3] != None else False for i in batch]),
11         device=self.DEVICE, dtype=torch.bool)
12
13     non_final_states = torch.cat([i[3] for i in batch if i[3] != None])
14
15     # Prepare batch data
16     batch = list(zip(*batch))
17     state_batch = torch.cat(batch[0])    # Current states
18     action_batch = torch.cat(batch[1])   # Actions taken
19     reward_batch = torch.cat(batch[2])   # Rewards received
20
21     # Get predicted Q-values using policy network
22     predicted_value = self.Policy_net(state_batch).gather(1, action_batch)
23
24     temp = torch.zeros(self.BATCH_SIZE).to(self.DEVICE)
25
26     # Calculate target Q-values using target network
27     with torch.no_grad():
28         temp[non_final_index] = self.Target_net(non_final_states).max(1)[0]
29
30     # Compute target values using Bellman equation
31     target_value = reward_batch + self.GAMMA * temp
32
33     # Calculate loss
34     criteria = torch.nn.SmoothL1Loss()
35     loss = criteria(predicted_value, target_value.unsqueeze(1))
36
37     # Update the policy network
38     self.optimizer.zero_grad()
39     loss.backward()
40     self.optimizer.step()
```

The training process implements the core DQN algorithm. The following are key points:

- Samples random batches from experience replay.
- Use policy network to get predicted value.
- Uses target network to get initial target values
- Apply Bellman equation on initial target values to compute target values
- Calculate loss based on target values and predicted values
- Update the policy model based on this loss.

2.5.1 Training Loop Implementation

```

1 def loop(self):
2     max_ep_len = 500 # maximum length of episode
3     max_training_timesteps = int(1e5) # maximum number of training timesteps
4     save_freq = 400*2 # save frequency
5     log_running_reward = 0
6     log_running_episodes = 0
7
8     while self.steps <= max_training_timesteps:
9         state, _ = self.env.reset()
10        state = torch.tensor(state, dtype=torch.float32,
11                             device=self.DEVICE).unsqueeze(0)
12
13        total_reward = 0
14
15        for t in range(1, max_ep_len+1):
16            action = self.choose_action(state) # get action
17            observation, reward, terminated, truncated, _ = \
18                self.env.step(action.item())
19            done = terminated or truncated # execute action to get next state and
20                reward
21            total_reward += reward
22
23            reward = torch.tensor([reward], device=self.DEVICE)
24            if terminated:
25                next_state = None
26            else:
27                next_state = torch.tensor(observation, dtype=torch.float32,
28                                         device=self.DEVICE).unsqueeze(0)
29
30            self.buffer.add(state, action, reward, next_state) # add the (state,
31                action, reward, next_state)
32            state = next_state
33            self.train() # start training model
34            self.soft_update() # update the model
35
36            if self.steps % save_freq == 0: #Save the average rewards (For
37                plotting purposes)
38                log_avg_reward = log_running_reward / log_running_episodes
39                log_avg_reward = round(log_avg_reward, 4)
40                self.avg_reward.append((self.steps, log_avg_reward))
41                log_running_reward = 0
42                log_running_episodes = 0
43
44            if done: # if episode is terminated we dont need to go to max_ep_len
45                break
46
47        log_running_reward += total_reward
48        log_running_episodes += 1

```

The training loop implements the core interaction between the agent and environment, managing both episode execution and logging. Key components include:

- **Training Parameters:**
 - Maximum episode length: 500 steps
 - Total training timesteps: 100,000
 - Logging frequency: Every 800 steps
- **Episode Structure:**

- Resets environment at start of each episode
- Converts state observations to tensors
- Tracks cumulative reward per episode

- **Step-by-Step Execution:**

- Action selection using epsilon-greedy policy
- Environment interaction and reward collection
- State processing and storage in replay buffer
- Training update and target network soft update

- **Performance Logging:**

- Maintains running average of episode rewards
- Logs performance metrics at regular intervals
- Handles both natural termination and truncation

2.6 Visualization and Analysis

```

1 def plot(self, name):
2     import matplotlib.pyplot as plt
3     plt.xlabel("Time steps")
4     plt.ylabel("Avg episodic reward")
5     plt.title(f"DQN plot for {name}")
6     plt.plot([i[0] for i in self.avg_reward],
7              [i[1] for i in self.avg_reward])
8     plt.show()

```

The visualization component has following important points:

- Plots average episodic reward against training timesteps
- Helps identify learning stability and convergence
- Facilitates comparison between different hyperparameter settings

2.7 Usage Example

```

1 env = gym.make("CartPole-v1")
2 cart = DQN(env, Policy)
3 cart.loop()
4 cart.plot("CartPole-V1")

```

The implementation demonstrates practical usage with the CartPole-v1 environment:

- Environment initialization using OpenAI Gym
- DQN agent creation with custom policy network
- Training execution through the main loop
- Performance visualization using matplotlib