

# Reinforcement Learning Report

Dokeun Lee

## Problem 1 - DQN Implementation

(a) Complete the DQN skeleton code.

- QNetwork

```
class QNetwork(nn.Module):
    def __init__(self, state_dim, action_dim):
        super(QNetwork, self).__init__()
        self.fc1 = nn.Linear(state_dim, 256)
        self.fc_out = nn.Linear(256, action_dim)
        self.relu = nn.ReLU()

    def forward(self, s):
        s = self.relu(self.fc1(s))
        q = self.fc_out(s)

        return q
```

This is a neural network that takes a state as input and outputs the Q-value for each action. It uses one hidden layer with 256 nodes and a ReLU activation function.

- DQNAgent

```
class DQNAgent:
    def __init__(self, state_size, action_size, device):
        self.state_size = state_size
        self.action_size = action_size
        self.device = device

        # Do not modify these hyper-parameters
        self.epochs = 1000
        self.discount_factor = 0.98
        self.learning_rate = 0.001 # learning rate for q function
        self.epsilon = 1.0 # initial epsilon value
        self.epsilon_min = 0.001 # minimum epsilon value
        self.batch_size = 256
        self.train_start = self.batch_size * 5
        self.memory = deque(maxlen=100000) # replay memory

        # You can modify this depending on environments.
        self.epsilon_decay_rate = 0.995 # decay rate

        # Define and initialize your networks and optimizer
        self.q_net = QNetwork(state_size, action_size).to(device)
        self.q_target = QNetwork(state_size, action_size).to(device)
        self.optimizer = optim.Adam(self.q_net.parameters(), lr=self.learning_rate)
        self.loss_function = nn.MSELoss()

    def update_target_network(self):
        # implement target q network update function
        self.q_target.load_state_dict(self.q_net.state_dict())

    def get_action(self, state, use_epsilon_greedy=True):
        if use_epsilon_greedy and random.random() < self.epsilon:
            # implement epsilon greedy policy given state
            action = random.randrange(self.action_size)
        else:
            # implement greedy policy given state
            # this greedy policy is used for evaluation
            state_tensor = torch.FloatTensor(state).unsqueeze(0).to(self.device)
            q_values = self.q_net(state_tensor)
            action = q_values.argmax().item()

        return action

    def append_sample(self, state, action, reward, next_state, done):
        # implement storing function given (s,a,r,s',done) into the replay memory.
        self.memory.append((state, action, reward, next_state, done))

    def get_samples(self, n):
        # implement transition random sampling function from the replay memory.
        # and make the transition to batch.
        samples = random.sample(self.memory, n)
        s_batch = torch.tensor([s[0] for s in samples]).float().to(self.device)
        a_batch = torch.tensor([s[1] for s in samples]).long().to(self.device)
        r_batch = torch.tensor([s[2] for s in samples]).float().to(self.device)
        s_next_batch = torch.tensor([s[3] for s in samples]).float().to(self.device)
        done_batch = torch.tensor([s[4] for s in samples]).float().to(self.device)
        # s,a,r,s',done (batch size, state dim)
        return s_batch, a_batch, r_batch, s_next_batch, done_batch

    def epsilon_decay(self):
        # implement epsilon decaying function
        self.epsilon = max(self.epsilon * self.epsilon_decay_rate, self.epsilon_min)

    def train(self):
        if len(self.memory) < self.train_start:
            return
        s_batch, a_batch, r_batch, s_next_batch, done_batch = self.get_samples(self.batch_size)

        # implement DQN training function.
        # You can return any statistics you want to check and analyze the training. (i.e. loss, q values, target q values, ... etc)
        target_values = r_batch + self.discount_factor * self.q_target(s_next_batch).detach().max(1)[0].unsqueeze(1) * (1 - done_batch)
        current_values = self.q_net(s_batch).gather(1, a_batch)
        loss = self.loss_function(current_values, target_values)

        self.optimizer.zero_grad()
        loss.backward()
        self.optimizer.step()
```

**Replay Memory:** It uses a replay buffer (**memory**) to store transition tuples of (**state**, **action**, **reward**, **next\_state**, **done**).

**Target Network:** To ensure training stability, a separate target network (**q\_target**) is used in addition to the main network (**q\_net**).

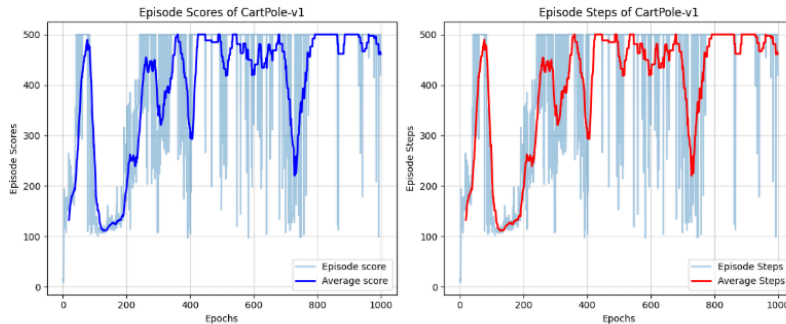
**get\_action:** It balances exploration and exploitation using an epsilon-greedy policy.

**train:** The agent learns by sampling a minibatch from the replay buffer. The loss is calculated as the Mean Squared Error (MSE) between the TD target computed by the target network and the current Q-value predicted by the main network.

(b) Train the DQN Agent on the three environments.

(c) Plot evaluation episode score and steps graphs on each environment.

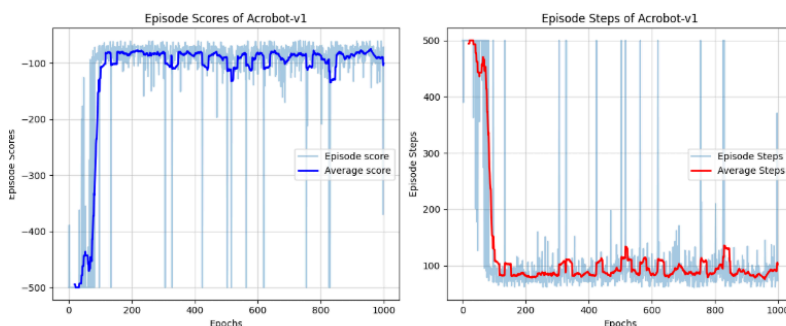
- Cartpole-v1, `self.epsilon_decay_rate = 0.995`



```
state size:4, action size:2, max_episode_steps:500, device:cuda
Epoch:0, Episode score:16.0, Episode steps:16, epsilon:0.995
Epoch:10, Episode score:107.0, Episode steps:107, epsilon:0.946354579813443
Epoch:20, Episode score:147.0, Episode steps:147, epsilon:0.908874278732445
Epoch:30, Episode score:178.0, Episode steps:178, epsilon:0.856882709551227
Epoch:40, Episode score:321.0, Episode steps:321, epsilon:0.8142285204175609
Epoch:50, Episode score:500.0, Episode steps:500, epsilon:0.774420942823988
Epoch:60, Episode score:293.0, Episode steps:293, epsilon:0.73559652608221
Epoch:70, Episode score:500.0, Episode steps:500, epsilon:0.7005493475733617
Epoch:80, Episode score:500.0, Episode steps:500, epsilon:0.6662995813682115
Epoch:90, Episode score:112.0, Episode steps:112, epsilon:0.6337242817644885
Epoch:100, Episode score:192.0, Episode steps:192, epsilon:0.6027415843882742
Epoch:110, Episode score:181.0, Episode steps:181, epsilon:0.5732736268885887
Epoch:120, Episode score:114.0, Episode steps:114, epsilon:0.5452463546625918
Epoch:130, Episode score:116.0, Episode steps:116, epsilon:0.5158933094845482
Epoch:140, Episode score:121.0, Episode steps:121, epsilon:0.4932355662165453
Epoch:150, Episode score:136.0, Episode steps:136, epsilon:0.4691213473457726
Epoch:160, Episode score:122.0, Episode steps:122, epsilon:0.44818065444672
Epoch:170, Episode score:138.0, Episode steps:138, epsilon:0.4243708486280985
Epoch:180, Episode score:123.0, Episode steps:123, epsilon:0.4036245827390106
Epoch:190, Episode score:148.0, Episode steps:148, epsilon:0.3838914347919885
Epoch:200, Episode score:171.0, Episode steps:171, epsilon:0.3651230261753626
Epoch:210, Episode score:252.0, Episode steps:252, epsilon:0.3472722151892322
Epoch:220, Episode score:228.0, Episode steps:228, epsilon:0.3302941218954743
Epoch:960, Episode score:500.0, Episode steps:500, epsilon:0.00809597512943647
Epoch:970, Episode score:500.0, Episode steps:500, epsilon:0.007695940256818795
Epoch:980, Episode score:500.0, Episode steps:500, epsilon:0.007318839301892576
Epoch:990, Episode score:500.0, Episode steps:500, epsilon:0.006961022023218655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

The agent learned quickly and succeeded in keeping the episode score close to 500. The training is not perfectly stable, but it attempts to maintain its performance.

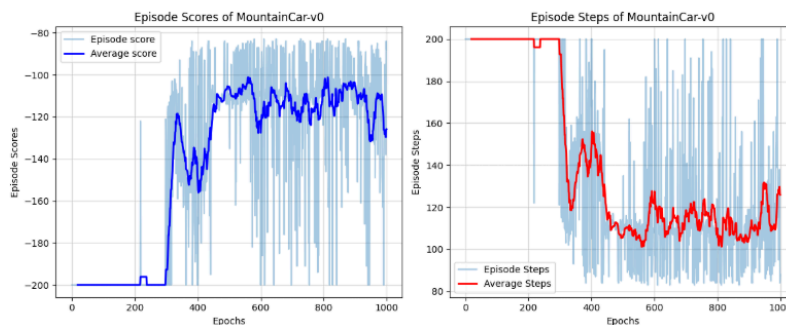
- Acrobot-v1, `self.epsilon_decay_rate = 0.995`



```
state size:8, action size:3, max_episode_steps:500, device:cuda
Epoch:0, Episode score:-500.0, Episode steps:500, epsilon:0.995
Epoch:10, Episode score:-500.0, Episode steps:500, epsilon:0.946354579813443
Epoch:20, Episode score:-500.0, Episode steps:500, epsilon:0.908874278732445
Epoch:30, Episode score:-500.0, Episode steps:500, epsilon:0.856882709551227
Epoch:40, Episode score:-500.0, Episode steps:500, epsilon:0.8142285204175609
Epoch:50, Episode score:-500.0, Episode steps:500, epsilon:0.774420942823988
Epoch:60, Episode score:-500.0, Episode steps:500, epsilon:0.73559652608221
Epoch:70, Episode score:-500.0, Episode steps:500, epsilon:0.7005493475733617
Epoch:80, Episode score:-106.0, Episode steps:107, epsilon:0.6662995813682115
Epoch:90, Episode score:92.0, Episode steps:93, epsilon:0.6337242817644885
Epoch:100, Episode score:88.0, Episode steps:89, epsilon:0.6027415843882742
Epoch:110, Episode score:76.0, Episode steps:77, epsilon:0.5732736268885887
Epoch:120, Episode score:-78.0, Episode steps:71, epsilon:0.5452463546625918
Epoch:130, Episode score:-78.0, Episode steps:79, epsilon:0.5158933094845482
Epoch:140, Episode score:-91.0, Episode steps:82, epsilon:0.4932355662165453
Epoch:150, Episode score:-99.0, Episode steps:100, epsilon:0.4691213473457726
Epoch:160, Episode score:-71.0, Episode steps:72, epsilon:0.44818065444672
Epoch:170, Episode score:-83.0, Episode steps:84, epsilon:0.4243708486280985
Epoch:180, Episode score:-99.0, Episode steps:100, epsilon:0.4036245827390106
Epoch:190, Episode score:-148.0, Episode steps:109, epsilon:0.3838914347919885
Epoch:200, Episode score:-69.0, Episode steps:79, epsilon:0.3651230261753626
Epoch:210, Episode score:-82.0, Episode steps:83, epsilon:0.3472722151892322
Epoch:220, Episode score:145.0, Episode steps:106, epsilon:0.3302941218954743
Epoch:960, Episode score:-188.0, Episode steps:101, epsilon:0.00809597512943647
Epoch:970, Episode score:-98.0, Episode steps:99, epsilon:0.007695940256818795
Epoch:980, Episode score:-61.0, Episode steps:62, epsilon:0.007318839301892576
Epoch:990, Episode score:-94.0, Episode steps:85, epsilon:0.006961022023218655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

While there is some variance, the agent tries to maintain a score of -100, and the average reward converged to around this level as training progressed.

- MountainCar-v0, `self.epsilon_decay_rate = 0.995`



```
state size:2, action size:3, max_episode_steps:200, device:cuda
Epoch:0, Episode score:-200.0, Episode steps:200, epsilon:0.995
Epoch:10, Episode score:-200.0, Episode steps:200, epsilon:0.946354579813443
Epoch:20, Episode score:-200.0, Episode steps:200, epsilon:0.908874278732445
Epoch:30, Episode score:-200.0, Episode steps:200, epsilon:0.856882709551227
Epoch:40, Episode score:-200.0, Episode steps:200, epsilon:0.8142285204175609
Epoch:50, Episode score:-200.0, Episode steps:200, epsilon:0.774420942823988
Epoch:60, Episode score:-200.0, Episode steps:200, epsilon:0.73559652608221
Epoch:70, Episode score:-200.0, Episode steps:200, epsilon:0.7005493475733617
Epoch:80, Episode score:-200.0, Episode steps:200, epsilon:0.6662995813682115
Epoch:90, Episode score:-200.0, Episode steps:200, epsilon:0.6337242817644885
Epoch:100, Episode score:-200.0, Episode steps:200, epsilon:0.6027415843882742
Epoch:110, Episode score:-200.0, Episode steps:200, epsilon:0.5732736268885887
Epoch:120, Episode score:-200.0, Episode steps:200, epsilon:0.5452463546625918
Epoch:130, Episode score:-200.0, Episode steps:200, epsilon:0.5158933094845482
Epoch:140, Episode score:-200.0, Episode steps:200, epsilon:0.4932355662165453
Epoch:150, Episode score:-200.0, Episode steps:200, epsilon:0.4691213473457726
Epoch:160, Episode score:-200.0, Episode steps:200, epsilon:0.44818065444672
Epoch:170, Episode score:-200.0, Episode steps:200, epsilon:0.4243708486280985
Epoch:180, Episode score:-200.0, Episode steps:200, epsilon:0.4036245827390106
Epoch:190, Episode score:-200.0, Episode steps:200, epsilon:0.3838914347919885
Epoch:200, Episode score:-200.0, Episode steps:200, epsilon:0.3651230261753626
Epoch:210, Episode score:-200.0, Episode steps:200, epsilon:0.3472722151892322
Epoch:220, Episode score:-200.0, Episode steps:200, epsilon:0.3302941218954743
Epoch:960, Episode score:-107.0, Episode steps:107, epsilon:0.00809597512943647
Epoch:970, Episode score:-85.0, Episode steps:85, epsilon:0.007695940256818795
Epoch:980, Episode score:-122.0, Episode steps:122, epsilon:0.007318839301892576
Epoch:990, Episode score:-200.0, Episode steps:200, epsilon:0.006961022023218655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

Starting from around 300 epochs, the performance gradually started to improve, converging to an average reward of around -120 with a maximum reward of -100.

(d) Are the DQN agents trained well? Report and analyze the experiment's results.

The DQN agent was trained successfully on the CartPole-v1 environment. Furthermore, it demonstrated meaningful learning on both Acrobot-v1 and MountainCar-v0, achieving convergence to scores around -100 and -120, respectively. However, the variance during training and the failure to reach a perfectly optimal score can be attributed to DQN's inherent problem of Overestimation Bias. This tendency to optimistically overestimate action values hinders training stability and limits the ability to find the optimal policy in complex or sparse-reward environments.

## Problem 2 - Double DQN Implementation

(a) Complete the Double DQN skeleton code.

- QNetwork

```
class QNetwork(nn.Module):
    def __init__(self, state_dim, action_dim):
        super(QNetwork, self).__init__()
        self.fc1 = nn.Linear(state_dim, 256)
        self.fc_out = nn.Linear(256, action_dim)
        self.relu = nn.ReLU()

    def forward(self, s):
        s = self.relu(self.fc1(s))
        q = self.fc_out(s)

        return q
```

- DoubleDQNAgent

```
class DoubleDQNAgent:
    def __init__(self, state_size, action_size, device):
        self.state_size = state_size
        self.action_size = action_size
        self.device = device

        # Do not modify these hyper-parameters
        self.Epochs = 1000
        self.discount_factor = 0.98
        self.learning_rate = 0.001 # learning rate for q function
        self.epsilon = 1.0 # initial epsilon value
        self.epsilon_min = 0.001 # minimum epsilon value
        self.batch_size = 256
        self.train_start = self.batch_size * 5
        self.memory = deque(maxlen=100000) # replay memory

        # You can modify this depending on environments.
        self.epsilon_decay_rate = 0.995 # decay rate

        # Define and initialize your networks and optimizer
        self.q_net = QNetwork(state_size, action_size).to(device)
        self.q_target = QNetwork(state_size, action_size).to(device)
        self.optimizer = optim.Adam(self.q_net.parameters(), lr=self.learning_rate)
        self.loss_function = nn.MSELoss()

    def update_target_network(self):
        # implement target Q network update function
        self.q_target.load_state_dict(self.q_net.state_dict())

    def get_action(self, state, use_epsilon_greedy=True):
        if use_epsilon_greedy and random.random() < self.epsilon:
            # Implement epsilon greedy policy given state
            action = random.randrange(self.action_size)
        else:
            # Implement greedy policy given state
            # this greedy policy is used for evaluation
            state_tensor = torch.FloatTensor(state).unsqueeze(0).to(self.device)
            q_values = self.q_net(state_tensor)
            action = q_values.argmax().item()

        return action

    def append_sample(self, state, action, reward, next_state, done):
        # implement storing function given (s,a,r,s',done) into the replay memory.
        self.memory.append((state, action, reward, next_state, done))

    def get_samples(self, n):
        # implement transition random sampling function from the replay memory,
        # and make the transition to batch.
        samples = random.sample(self.memory, n)
        s_batch = torch.tensor([s[0] for s in samples]).float().to(self.device)
        a_batch = torch.tensor([s[1] for s in samples]).long().to(self.device)
        r_batch = torch.tensor([s[2] for s in samples]).float().to(self.device)
        s_next_batch = torch.tensor([s[3] for s in samples]).float().to(self.device)
        done_batch = torch.tensor([s[4] for s in samples]).float().to(self.device)
        # i.e.) s_batch : (batch_size, state_dim)
        return s_batch, a_batch, r_batch, s_next_batch, done_batch

    def epsilon_decay(self):
        # implement epsilon decaying function
        self.epsilon = max(self.epsilon * self.epsilon_decay_rate, self.epsilon_min)

    def train(self):
        if len(self.memory) < self.train_start:
            return
        s_batch, a_batch, r_batch, s_next_batch, done_batch = self.get_samples(self.batch_size)

        best_actions = self.q_net(s_next_batch).argmax(dim=1).unsqueeze(1)
        next_q_values = self.q_target(s_next_batch).gather(1, best_actions).detach()
        target_values = r_batch + self.discount_factor * next_q_values * (1 - done_batch)

        current_values = self.q_net(s_batch).gather(1, a_batch)
        loss = self.loss_function(current_values, target_values)

        self.optimizer.zero_grad()
        loss.backward()
        self.optimizer.step()
```

To address DQN's overestimation bias, a Double DQN agent was implemented.

- The network architecture is identical to DQN, using two networks: a main network (`q_net`) and a target network (`q_target`).
- The key difference lies in the TD target calculation within the `train` method. Unlike DQN, the role of selecting the best next action is handled by the main network, while the role of evaluating the value of that selected action is handled by the target network.

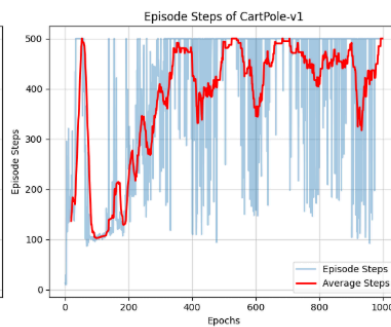
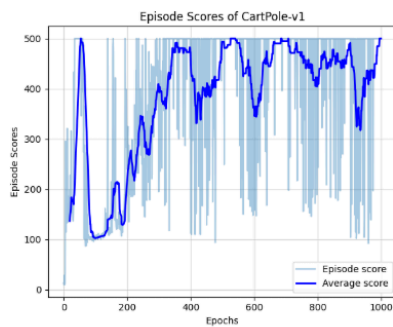
$$Q_1(s,a) \leftarrow Q_1(s,a) + \alpha [r + \gamma Q_2(s', \operatorname{argmax}_{a'} Q_1(s', a')) - Q_1(s,a)]$$

$$Q_2(s,a) \leftarrow Q_2(s,a) + \alpha [r + \gamma Q_1(s', \operatorname{argmax}_{a'} Q_2(s', a')) - Q_2(s,a)]$$

(b) Train the Double DQN Agent on the three environments.

(c) Plot evaluation episode score and steps graphs on each environment.

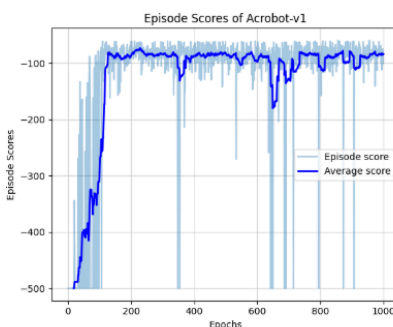
- Cartpole-v1, `self.epsilon_decay_rate = 0.995`



```
state size:4, action size:2, max episode steps:500, device:cuda
Epoch:0, Episode score:13.0, Episode steps:13, epsilon:0.995
Epoch:10, Episode score:13.0, Episode steps:13, epsilon:0.946354579813443
Epoch:20, Episode score:228.0, Episode steps:228, epsilon:0.9088874278732445
Epoch:30, Episode score:154.0, Episode steps:154, epsilon:0.856882769551227
Epoch:40, Episode score:500.0, Episode steps:500, epsilon:0.8142282584175669
Epoch:50, Episode score:500.0, Episode steps:500, epsilon:0.7744209942832888
Epoch:60, Episode score:394.0, Episode steps:394, epsilon:0.73655952988221
Epoch:70, Episode score:164.0, Episode steps:164, epsilon:0.7004849472730617
Epoch:80, Episode score:106.0, Episode steps:106, epsilon:0.6662959513682115
Epoch:90, Episode score:183.0, Episode steps:183, epsilon:0.633742817644086
Epoch:100, Episode score:102.0, Episode steps:102, epsilon:0.6027415843082742
Epoch:110, Episode score:106.0, Episode steps:106, epsilon:0.5732736268885887
Epoch:120, Episode score:108.0, Episode steps:108, epsilon:0.5452463546625918
Epoch:130, Episode score:162.0, Episode steps:162, epsilon:0.5185933094454582
Epoch:140, Episode score:118.0, Episode steps:118, epsilon:0.493235662105453
Epoch:150, Episode score:115.0, Episode steps:115, epsilon:0.4691213473457726
Epoch:160, Episode score:179.0, Episode steps:179, epsilon:0.446186862443672
Epoch:170, Episode score:122.0, Episode steps:122, epsilon:0.4243728466268885
Epoch:180, Episode score:138.0, Episode steps:138, epsilon:0.4036245882390186
Epoch:190, Episode score:162.0, Episode steps:162, epsilon:0.38389143477919885
Epoch:200, Episode score:264.0, Episode steps:264, epsilon:0.3651283261753626
Epoch:210, Episode score:265.0, Episode steps:265, epsilon:0.347272151889232
Epoch:220, Episode score:210.0, Episode steps:210, epsilon:0.3302941218954743
Epoch:300, Episode score:500.0, Episode steps:500, epsilon:0.688995957512943647
Epoch:350, Episode score:500.0, Episode steps:500, epsilon:0.607695042568818795
Epoch:400, Episode score:500.0, Episode steps:500, epsilon:0.607318839301892576
Epoch:450, Episode score:500.0, Episode steps:500, epsilon:0.606961922012318655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

Similar to DQN, the agent learned quickly and succeeded in keeping the episode score close to 500. The training is not perfectly stable, but it attempts to maintain its performance.

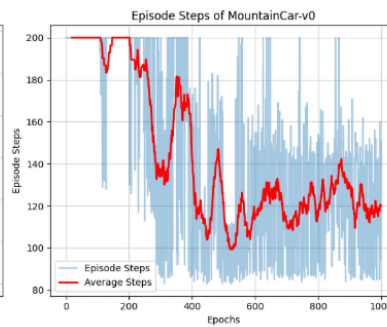
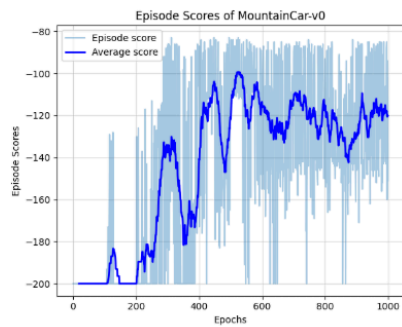
- Acrobot-v1, `self.epsilon_decay_rate = 0.995`



```
state size:6, action size:3, max episode steps:500, device:cuda
Epoch:0, Episode score:-500.0, Episode steps:500, epsilon:0.995
Epoch:10, Episode score:-500.0, Episode steps:500, epsilon:0.946354579813443
Epoch:20, Episode score:-500.0, Episode steps:500, epsilon:0.9088874278732445
Epoch:30, Episode score:-500.0, Episode steps:500, epsilon:0.856882769551227
Epoch:40, Episode score:-500.0, Episode steps:500, epsilon:0.8142282584175669
Epoch:50, Episode score:-500.0, Episode steps:500, epsilon:0.7744209942832888
Epoch:60, Episode score:-500.0, Episode steps:500, epsilon:0.73655952988221
Epoch:70, Episode score:-91.0, Episode steps:92, epsilon:0.7004849472730617
Epoch:80, Episode score:-182.0, Episode steps:163, epsilon:0.6662959513682115
Epoch:90, Episode score:-121.0, Episode steps:194, epsilon:0.633742817644086
Epoch:100, Episode score:-183.0, Episode steps:184, epsilon:0.6027415843082742
Epoch:110, Episode score:-81.0, Episode steps:82, epsilon:0.5732736268885887
Epoch:120, Episode score:-83.0, Episode steps:84, epsilon:0.5452463546625918
Epoch:130, Episode score:-86.0, Episode steps:87, epsilon:0.5185933094454582
Epoch:140, Episode score:-115.0, Episode steps:115, epsilon:0.493235662105453
Epoch:150, Episode score:-84.0, Episode steps:95, epsilon:0.4691213473457726
Epoch:160, Episode score:-82.0, Episode steps:83, epsilon:0.446186862443672
Epoch:170, Episode score:-89.0, Episode steps:76, epsilon:0.4243728466268885
Epoch:180, Episode score:-114.0, Episode steps:115, epsilon:0.4036245882390186
Epoch:190, Episode score:-92.0, Episode steps:93, epsilon:0.38389143477919885
Epoch:200, Episode score:-61.0, Episode steps:62, epsilon:0.3651283261753626
Epoch:210, Episode score:-77.0, Episode steps:78, epsilon:0.347272151889232
Epoch:220, Episode score:-82.0, Episode steps:83, epsilon:0.3302941218954743
Epoch:300, Episode score:-90.0, Episode steps:91, epsilon:0.688995957512943647
Epoch:350, Episode score:-91.0, Episode steps:92, epsilon:0.607695042568818795
Epoch:400, Episode score:-86.0, Episode steps:87, epsilon:0.607318839301892576
Epoch:450, Episode score:-87.0, Episode steps:88, epsilon:0.606961922012318655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

A performance improvement over DQN is visible. The variance is smaller, indicating more stable training, although the maximum reward level is similar.

- MountainCar-v0, self.epsilon\_decay\_rate = 0.995



```
state size:2, action size:3, max episode steps:200, device:cuda
Epoch:0, Episode score:-200.0, Episode steps:200, epsilon:0.995
Epoch:10, Episode score:-200.0, Episode steps:200, epsilon:0.946354578811443
Epoch:20, Episode score:-200.0, Episode steps:200, epsilon:0.9098874278732445
Epoch:30, Episode score:-200.0, Episode steps:200, epsilon:0.8504822709951227
Epoch:40, Episode score:-200.0, Episode steps:200, epsilon:0.8142285284175609
Epoch:50, Episode score:-200.0, Episode steps:200, epsilon:0.7744209942832988
Epoch:60, Episode score:-200.0, Episode steps:200, epsilon:0.72655952608221
Epoch:70, Episode score:-200.0, Episode steps:200, epsilon:0.7005493479731617
Epoch:80, Episode score:-200.0, Episode steps:200, epsilon:0.666295513682115
Epoch:90, Episode score:-200.0, Episode steps:200, epsilon:0.633724281764488
Epoch:100, Episode score:-200.0, Episode steps:200, epsilon:0.602715543882742
Epoch:110, Episode score:-200.0, Episode steps:200, epsilon:0.573273626885887
Epoch:120, Episode score:-200.0, Episode steps:200, epsilon:0.5452463540625918
Epoch:130, Episode score:-200.0, Episode steps:200, epsilon:0.518393304845452
Epoch:140, Episode score:-200.0, Episode steps:200, epsilon:0.493235562625453
Epoch:150, Episode score:-200.0, Episode steps:200, epsilon:0.46912134373457726
Epoch:160, Episode score:-200.0, Episode steps:200, epsilon:0.44618662443872
Epoch:170, Episode score:-200.0, Episode steps:200, epsilon:0.4243720840620985
Epoch:180, Episode score:-200.0, Episode steps:200, epsilon:0.403624582389186
Epoch:190, Episode score:-200.0, Episode steps:200, epsilon:0.383834577315885
Epoch:200, Episode score:-200.0, Episode steps:200, epsilon:0.3651203261753626
Epoch:210, Episode score:-200.0, Episode steps:200, epsilon:0.3472722151889232
Epoch:220, Episode score:-200.0, Episode steps:200, epsilon:0.3302941211854743
...
Epoch:960, Episode score:-89.0, Episode steps:89, epsilon:0.608998597512943647
Epoch:970, Episode score:-149.0, Episode steps:149, epsilon:0.607859649254981795
Epoch:980, Episode score:-149.0, Episode steps:149, epsilon:0.607318839301892576
Epoch:990, Episode score:-189.0, Episode steps:189, epsilon:0.6069610222032818655
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

Starting from around 200 epochs, the performance gradually increased, which is faster than DQN. The maximum reward is also better, reaching around -95.

- (d) Are the Double DQN agents trained well? Report and analyze the experiment's results.

The Double DQN agent demonstrated improved performance compared to DQN overall. In the Acrobot environment, training stability was improved with less variance, and in the MountainCar environment, it achieved a faster learning speed and a higher maximum reward. This is because Double DQN effectively mitigates the overestimation bias of DQN by decoupling action selection and value evaluation. As a result, more accurate and reliable Q-value learning became possible, leading to enhanced overall performance and stability.

### Problem 3 - Dueling Double DQN Implementation

- (a) You can implement any RL algorithm that could increase performance, but experiment on the three environments.

- DuelingQNetwork

```
class DuelingQNetwork(nn.Module):
    def __init__(self, state_dim, action_dim):
        super(DuelingQNetwork, self).__init__()

        self.feature_layer = nn.Sequential(
            nn.Linear(state_dim, 256),
            nn.ReLU(),
            nn.Linear(256, 128),
            nn.ReLU()
        )

        self.value_stream = nn.Sequential(
            nn.Linear(128, 64),
            nn.ReLU(),
            nn.Linear(64, 1)
        )

        self.advantage_stream = nn.Sequential(
            nn.Linear(128, 64),
            nn.ReLU(),
            nn.Linear(64, action_dim)
        )

    def forward(self, state):
        features = self.feature_layer(state)
        v = self.value_stream(features)
        a = self.advantage_stream(features)
        q_values = v + (a - a.mean(dim=-1, keepdim=True))

        return q_values
```



- DuelingDoubleDQNAgent

```

class DuelingDoubleDQNAgent:
    def __init__(self, state_size, action_size, device):
        self.state_size = state_size
        self.action_size = action_size
        self.device = device

        # Do not modify these hyper-parameters
        self.Epochs = 1000
        self.discount_factor = 0.98
        self.learning_rate = 0.001 # learning rate for q function
        self.epsilon = 1.0 # initial epsilon value
        self.epsilon_min = 0.001 # minimum epsilon value
        self.batch_size = 256
        self.train_start = self.batch_size * 5
        self.memory = deque(maxlen=100000) # replay memory

        # You can modify this depending on environments.
        self.epsilon_decay_rate = 0.995 # decay rate

        self.q_net = DuelingQNetwork(state_size, action_size).to(device)
        self.q_target = DuelingQNetwork(state_size, action_size).to(device)
        self.optimizer = optim.Adam(self.q_net.parameters(), lr=self.learning_rate)

        self.update_target_network()

    def update_target_network(self):
        self.q_target.load_state_dict(self.q_net.state_dict())

    def get_action(self, state, use_epsilon_greedy=True):
        if use_epsilon_greedy and random.random() < self.epsilon:
            return random.randrange(self.action_size)
        else:
            with torch.no_grad():
                state_tensor = torch.FloatTensor(state).unsqueeze(0).to(self.device)
                q_values = self.q_net(state_tensor)
                return q_values.argmax().item()

    def append_sample(self, state, action, reward, next_state, done):
        self.memory.append((state, action, reward, next_state, done))

    def epsilon_decay(self):
        self.epsilon = max(self.epsilon * self.epsilon_decay_rate, self.epsilon_min)

    def train(self):
        if len(self.memory) < self.train_start:
            return

        batch = random.sample(self.memory, self.batch_size)
        states, actions, rewards, next_states, dones = zip(*batch)

        states = torch.FloatTensor(np.array(states)).to(self.device)
        actions = torch.LongTensor(actions).unsqueeze(1).to(self.device)
        rewards = torch.FloatTensor(rewards).unsqueeze(1).to(self.device)
        next_states = torch.FloatTensor(np.array(next_states)).to(self.device)
        dones = torch.FloatTensor(dones).unsqueeze(1).to(self.device)

        with torch.no_grad():
            best_actions = self.q_net(next_states).argmax(dim=1).unsqueeze(1)
            next_q_values = self.q_target(next_states).gather(1, best_actions)
            target_q_values = rewards + (1 - dones) * self.discount_factor * next_q_values

        current_q_values = self.q_net(states).gather(1, actions)
        loss = nn.MSELoss(current_q_values, target_q_values)

        self.optimizer.zero_grad()
        loss.backward()
        self.optimizer.step()

```

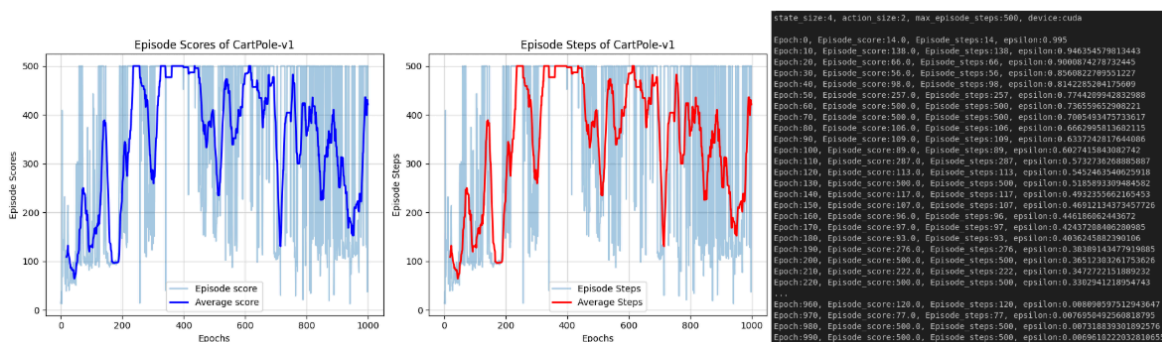
To further improve performance, a Dueling Double DQN agent was implemented by combining the Dueling Network architecture with Double DQN.

- **DuelingQNetwork:** This network is internally divided into two streams.
  - Value Stream: Calculates the value of the state itself,  $V(s)$ .
  - Advantage Stream: Calculates the advantage for each action over the average,  $A(s, a)$ .
- The final Q-value is calculated by combining these two streams using the formula:  $Q(s, a) = V(s) + (A(s, a) - \text{mean}(A(s, a)))$ . This allows for more efficient learning by decoupling the state's value from the actions.
- The rest of the agent follows the learning method of Double DQN.

(b) Train the Dueling Double DQN Agent on the three environments.

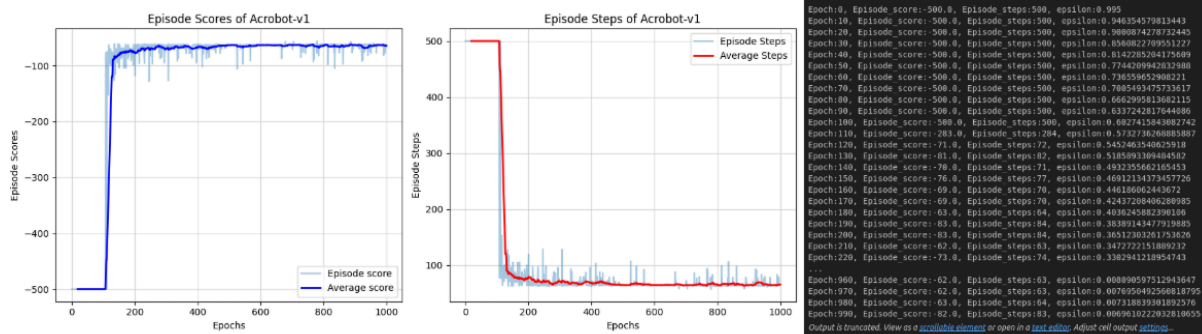
(c) Plot evaluation episode score and steps graphs on each environment.

- Cartpole-v1, self.epsilon\_decay\_rate = 0.995



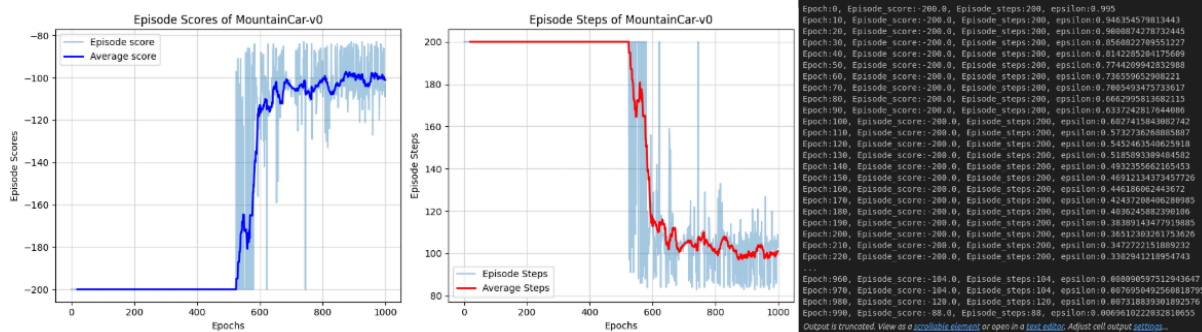
Similar to DQN and Double DQN, it converges quickly and attempts to maintain performance, but its performance is slightly worse than the other two algorithms.

- Acrobot-v1, `self.epsilon_decay_rate = 0.995`



Compared to the other two algorithms, the variance is definitively reduced. It consistently attempts to maintain a reward of -60, which seems to perfectly solve a certain drawback of DQN and Double DQN.

- MountainCar-v0, `self.epsilon_decay_rate = 0.995`



Its convergence speed is slower than DQN and Double DQN. Performance gradually stabilizes from 500 epochs, and it was confirmed to have excellent performance, attempting to maintain a superior average reward of -100 compared to the other two algorithms.

- (d) Are the Dueling Double DQN agents trained well? Report and analyze the experiment's results.

The Dueling Double DQN agent showed mixed but insightful results depending on the environment. In the Acrobot environment, it significantly reduced variance to stably maintain a high score (-60), and in the MountainCar environment, it ultimately achieved the best average score (-100) despite a slower start. This demonstrates the strength of the Dueling Network in finding more stable and qualitatively higher policies in complex problems by efficiently learning state values. Conversely, for the relatively simple CartPole environment, the added network complexity may have led to a slight degradation in performance. Nevertheless, the overwhelming stability and superior performance in the more difficult tasks can be attributed to the synergistic effect of Double DQN's bias reduction and the Dueling Network's efficient value learning.