

DA6400: Reinforcement Learning Programming Assignment 2 Report

KUSH RUSHABH SHAH | CH21B053 M S SAI TEJA | ME21B171

April 14, 2025

Abstract

This report describes the implementation and evaluation of two reinforcement learning algorithms on Gymnasium environments. In particular, we implemented two variants of Dueling-DQN (Type-1 and Type-2) and two variants of the Monte Carlo REINFORCE algorithm (with and without a TD(0) baseline) on Acrobot-v1 and CartPole-v1. The experimental results are averaged over 5 random seeds, and performance is compared in terms of episodic returns and variance.

1 Introduction

In this assignment, we are implementing two families of RL algorithms:

- **Dueling-DQN:** where the Q-value is decomposed into a state value function $V(s)$ and an advantage $A(s, a)$. Two update rules are implemented:

- **Type-1:**

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

- **Type-2:**

$$Q(s, a; \theta) = V(s; \theta) + \left(A(s, a; \theta) - \max_{a' \in A} A(s, a'; \theta) \right)$$

- **Monte-Carlo REINFORCE:** where the policy is updated using entire episode returns. Two variants are implemented:

- Without baseline:

$$\theta \leftarrow \theta + \alpha G_t \frac{\nabla \pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta)}$$

- With baseline, where the baseline $V(s; \Phi)$ is updated using TD(0):

$$\theta \leftarrow \theta + \alpha (G_t - V(s_t; \Phi)) \frac{\nabla \pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta)}$$

All experiments use a discount factor $\gamma = 0.99$ and are averaged over 5 random seeds.

2 Environments

We test the algorithms on:

Acrobot-v1: A two-link pendulum where the goal is to swing the free end above a given height.

CartPole-v1: A cart-pole system in which the objective is to balance an upright pole on a moving cart.

3 Algorithm Implementations

3.1 Dueling-DQN

The network architecture for Dueling-DQN uses a shared feature extraction layer and splits the output into:

1. A value stream $V(s)$.
2. An advantage stream $A(s,a)$.

The final Q-values are computed according to either Type-1 or Type-2 aggregation.

3.1.1 Dueling-DQN Code Snippet

```
1 class DuelingDQN(nn.Module):
2     def __init__(self, state_dim, action_dim, hidden_dim=128):
3         super(DuelingDQN, self).__init__()
4         % Shared feature extraction
5         self.feature = nn.Sequential(
6             nn.Linear(state_dim, hidden_dim),
7             nn.ReLU(),
8             nn.Linear(hidden_dim, hidden_dim),
9             nn.ReLU()
10        )
11        % Value stream: computes V(s)
12        self.value_stream = nn.Sequential(
13            nn.Linear(hidden_dim, hidden_dim//2),
14            nn.ReLU(),
15            nn.Linear(hidden_dim//2, 1)
16        )
17        % Advantage stream: computes A(s,a)
18        self.advantage_stream = nn.Sequential(
19            nn.Linear(hidden_dim, hidden_dim//2),
20            nn.ReLU(),
21            nn.Linear(hidden_dim//2, action_dim)
22        )
23
24    def forward(self, state, aggregation_type=1):
25        features = self.feature(state)
26        value = self.value_stream(features)
27        advantages = self.advantage_stream(features)
28        if aggregation_type == 1:
29            % Type-1: subtract the mean advantage
30            return value + (advantages - advantages.mean(dim=1, keepdim=True))
31        else:
32            % Type-2: subtract the maximum advantage
33            return value + (advantages - advantages.max(dim=1, keepdim=True)
34                           [0])
```

Listing 1: Dueling-DQN Network Architecture and Forward Pass

```
1 class DuelingDQNAgent:
2     """Agent implementing Dueling DQN with both Type-1 and Type-2 update rules.
3     """
4     def __init__(self, state_dim, action_dim, aggregation_type=1,
5                  lr=0.001, gamma=0.99, epsilon_start=1.0,
6                  epsilon_min=0.01, epsilon_decay=0.995,
7                  buffer_size=10000, batch_size=64,
8                  target_update=10, hidden_dim=128):
9         self.state_dim = state_dim
10        self.action_dim = action_dim
```

```

10     self.aggregation_type = aggregation_type
11     self.gamma = gamma
12     self.epsilon = epsilon_start
13     self.epsilon_min = epsilon_min
14     self.epsilon_decay = epsilon_decay
15     self.batch_size = batch_size
16     self.target_update = target_update
17     self.update_count = 0
18
19     # Initialize networks
20     self.policy_net = DuelingDQN(state_dim, action_dim, hidden_dim)
21     self.target_net = DuelingDQN(state_dim, action_dim, hidden_dim)
22     self.target_net.load_state_dict(self.policy_net.state_dict())
23     self.target_net.eval()
24
25     # Initialize optimizer
26     self.optimizer = optim.Adam(self.policy_net.parameters(), lr=lr)
27
28     # Initialize replay buffer
29     self.buffer = ReplayBuffer(buffer_size)
30
31     # Device configuration
32     self.device = torch.device("cuda" if torch.cuda.is_available() else "
        cpu")
33     self.policy_net.to(self.device)
34     self.target_net.to(self.device)
35
36     def select_action(self, state, training=True):
37         """Select action using epsilon-greedy policy."""
38         if training and np.random.rand() < self.epsilon:
39             # Exploration: random action
40             return np.random.randint(self.action_dim)
41         else:
42             # Exploitation: greedy action
43             state = torch.FloatTensor(state).unsqueeze(0).to(self.device)
44             with torch.no_grad():
45                 q_values = self.policy_net(state, self.aggregation_type)
46                 return q_values.argmax().item()
47
48     def update(self):
49         """Update the policy network using a batch of experiences."""
50         if len(self.buffer) < self.batch_size:
51             return 0
52
53         # Sample batch from replay buffer
54         states, actions, rewards, next_states, dones = self.buffer.sample(self.
            batch_size)
55
56         # Convert to tensors
57         states = torch.FloatTensor(states).to(self.device)
58         actions = torch.LongTensor(actions).to(self.device)
59         rewards = torch.FloatTensor(rewards).to(self.device)
60         next_states = torch.FloatTensor(next_states).to(self.device)
61         dones = torch.FloatTensor(dones).to(self.device)
62
63         # Get current Q values
64         q_values = self.policy_net(states, self.aggregation_type).gather(1,
            actions.unsqueeze(1))
65
66         # Get next Q values from target network
67         with torch.no_grad():
68             next_q_values = self.target_net(next_states, self.aggregation_type)
            .max(1)[0]

```

```

69
70     # Compute target Q values
71     target_q_values = rewards + (1 - dones) * self.gamma * next_q_values
72
73     # Compute loss
74     loss = F.mse_loss(q_values.squeeze(), target_q_values)
75
76     # Optimize
77     self.optimizer.zero_grad()
78     loss.backward()
79     self.optimizer.step()
80
81     # Update target network periodically
82     self.update_count += 1
83     if self.update_count % self.target_update == 0:
84         self.target_net.load_state_dict(self.policy_net.state_dict())
85
86     # Decay epsilon
87     self.epsilon = max(self.epsilon * self.epsilon_decay, self.epsilon_min)
88
89     return loss.item()
90
91 def train(self, env, num_episodes, max_steps=500):
92     """Train the agent over multiple episodes."""
93     episode_returns = []
94
95     for episode in range(num_episodes):
96         state, _ = env.reset()
97         episode_return = 0
98
99         for step in range(max_steps):
100             # Select and perform action
101             action = self.select_action(state)
102             next_state, reward, done, truncated, _ = env.step(action)
103
104             # Store transition in replay buffer
105             self.buffer.push(state, action, reward, next_state, done)
106
107             # Update state and episode return
108             state = next_state
109             episode_return += reward
110
111             # Update the network
112             loss = self.update()
113
114             if done or truncated:
115                 break
116
117             episode_returns.append(episode_return)
118             print(f"Episode {episode+1}, Return: {episode_return}, Epsilon: {
119                 self.epsilon:.4f}")
120
121     return episode_returns

```

Listing 2: Dueling-DQN Agent Code

3.2 Monte-Carlo REINFORCE

The REINFORCE agent collects complete episode returns, then updates the policy network using either:

- The full Monte Carlo return G_t (without baseline), or

- The advantage $G_t - V(s_t)$ (with baseline), where the value network is updated at each step using TD(0).

3.2.1 REINFORCE Agent Code Snippet

```

1 class REINFORCEAgent:
2     def __init__(self, state_dim, action_dim, use_baseline=False,
3                 lr_policy=1e-3, lr_value=1e-3, gamma=0.99, hidden_dim=128):
4         self.use_baseline = use_baseline
5         self.gamma = gamma
6         self.policy_net = PolicyNetwork(state_dim, action_dim, hidden_dim)
7         self.policy_optimizer = optim.Adam(self.policy_net.parameters(), lr=
            lr_policy)
8         if use_baseline:
9             self.value_net = ValueNetwork(state_dim, hidden_dim)
10            self.value_optimizer = optim.Adam(self.value_net.parameters(), lr=
                lr_value)
11        self.device = torch.device("cuda" if torch.cuda.is_available() else "
            cpu")
12        self.policy_net.to(self.device)
13        if use_baseline:
14            self.value_net.to(self.device)
15
16        % Update baseline using TD(0) at each step
17        def _update_baseline_td0(self, state, reward, next_state, done):
18            state_tensor = torch.FloatTensor(state).unsqueeze(0).to(self.device)
19            next_state_tensor = torch.FloatTensor(next_state).unsqueeze(0).to(self.
                device)
20            with torch.no_grad():
21                next_value = 0 if done else self.value_net(next_state_tensor)
22            reward_tensor = torch.tensor([reward], device=self.device, dtype=torch.
                float32)
23            td_target = reward_tensor + self.gamma * next_value
24            current_value = self.value_net(state_tensor)
25            value_loss = F.mse_loss(current_value, td_target)
26            self.value_optimizer.zero_grad()
27            value_loss.backward()
28            self.value_optimizer.step()
29
30        % Policy is updated after the episode using the complete trajectory
31        def _update_policy(self, states, rewards, log_probs):
32            returns = self.calculate_returns(rewards)
33            states = torch.FloatTensor(np.array(states)).to(self.device)
34            if self.use_baseline:
35                with torch.no_grad():
36                    state_values = self.value_net(states).squeeze(-1)
37                advantages = returns - state_values
38            else:
39                advantages = returns
40            if len(advantages) > 1:
41                advantages = (advantages - advantages.mean()) / (advantages.std() +
                    1e-8)
42            policy_loss = 0
43            for log_prob, advantage in zip(log_probs, advantages):
44                policy_loss -= log_prob * advantage
45            self.policy_optimizer.zero_grad()
46            policy_loss.backward()
47            self.policy_optimizer.step()
48
49        % Helper: compute Monte Carlo returns
50        def calculate_returns(self, rewards):
51            returns = []

```

```

52     R = 0
53     for r in reversed(rewards):
54         R = r + self.gamma * R
55         returns.insert(0, R)
56     return torch.tensor(returns, device=self.device, dtype=torch.float32)
57
58     def select_action(self, state):
59         state = torch.FloatTensor(state).to(self.device)
60         dist = self.policy_net(state)
61         action = dist.sample()
62         return action.item(), dist.log_prob(action)
63
64     def train(self, env, num_episodes, max_steps=500):
65         episode_returns = []
66         for episode in range(num_episodes):
67             state, _ = env.reset()
68             states, rewards, log_probs = [], [], []
69             episode_return = 0
70             for step in range(max_steps):
71                 action, log_prob = self.select_action(state)
72                 next_state, reward, done, truncated, _ = env.step(action)
73                 states.append(state)
74                 rewards.append(reward)
75                 log_probs.append(log_prob)
76                 if self.use_baseline:
77                     self._update_baseline_td0(state, reward, next_state, done
78                                                or truncated)
79                 state = next_state
80                 episode_return += reward
81                 if done or truncated:
82                     break
83             self._update_policy(states, rewards, log_probs)
84             episode_returns.append(episode_return)
85             print(f"Episode {episode+1}, Return: {episode_return}")
86         return episode_returns

```

Listing 3: REINFORCE with TD(0) Baseline

4 Experimental Setup

Experiments are run with $\gamma = 0.99$ and averaged over 5 random seeds to account for stochasticity. Hyperparameters are tuned by minimizing the regret in each experiment. The episodic return as a function of episode number is plotted, showing both the mean and variance across the seeds.

5 Results

5.1 Dueling-DQN Results

Figure 1 shows the mean episodic return of Type-1 and Type-2 variants on CartPole-v1, while Figure 2 shows the corresponding results on Acrobot-v1.

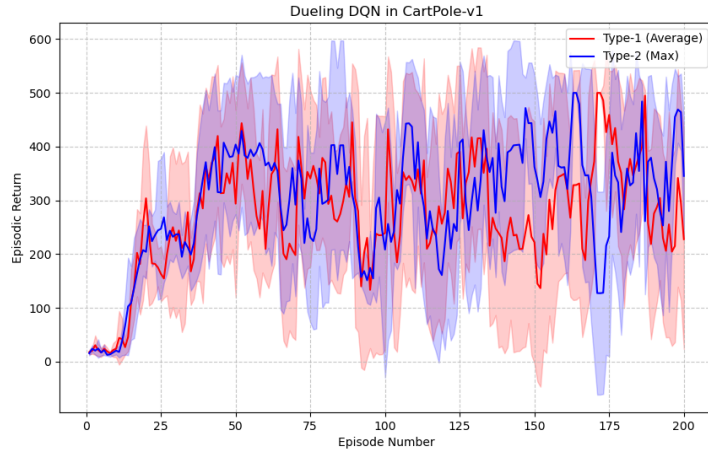


Figure 1: Comparison of Type-1 and Type-2 Dueling-DQN on CartPole-v1.

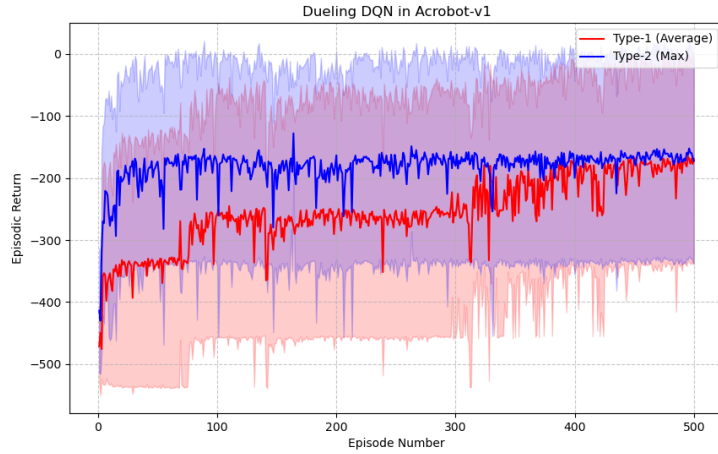


Figure 2: Comparison of Type-1 and Type-2 Dueling-DQN on Acrobot-v1.

5.2 Monte-Carlo REINFORCE Results

Figures 3 and 4 display the performance of REINFORCE with and without a baseline (using $TD(0)$) on CartPole-v1 and Acrobot-v1, respectively.

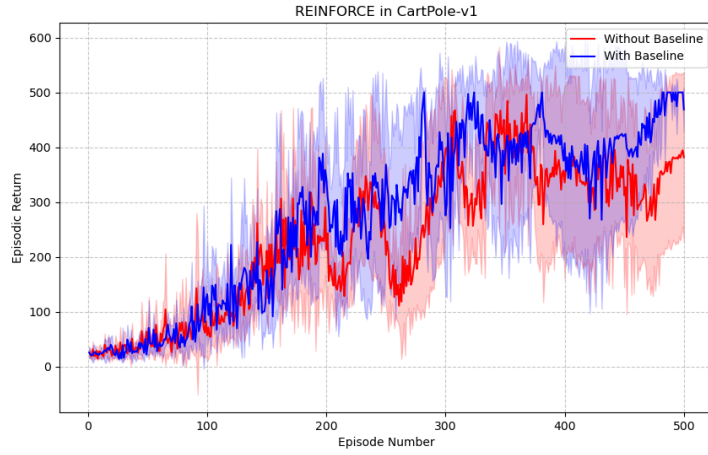


Figure 3: Monte-Carlo REINFORCE on CartPole-v1: With and Without Baseline.

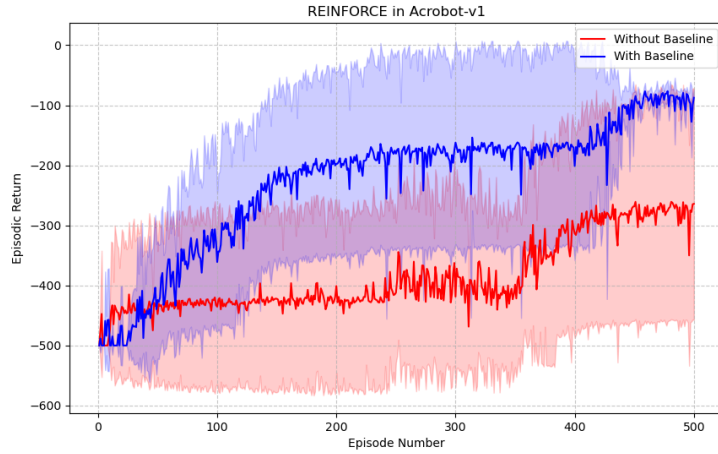


Figure 4: Monte-Carlo REINFORCE on Acrobot-v1: With and Without Baseline.

6 Hyperparameter Tuning

We tested various hyperparameter settings. The top setting for each algorithm were:

6.1 DQN Duelling - Cartpole-v1: Average

The hyperparameters used for DQN Duelling on Cartpole-v1 (Average) are:

- **Batch Size**
- **Epsilon decay**
- **Hidden Dimension**
- **Learning rate**
- **Target Update**

Figure 5 shows the parallel coordinates plot for the average performance.

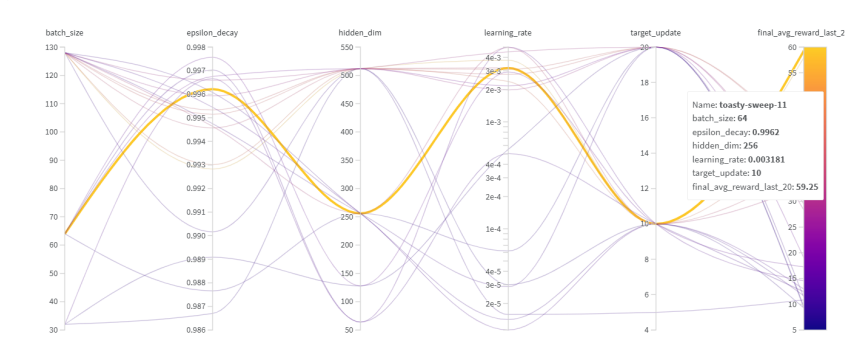


Figure 5: DQN Duelling - Cartpole-v1 (Average): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Batch Size:** 64
- **Epsilon decay:** 0.9962
- **Hidden Dimension:** 256
- **Learning rate:** 0.003181
- **Target Update:** 10

6.2 DQN Duelling - Cartpole-v1: Maximum

The hyperparameters used for DQN Duelling on Cartpole-v1 (Maximum) are:

- **Batch Size**
- **Epsilon decay**
- **Hidden Dimension**
- **Learning rate**
- **Target Update**

Figure 6 shows the parallel coordinates plot for the maximum performance.

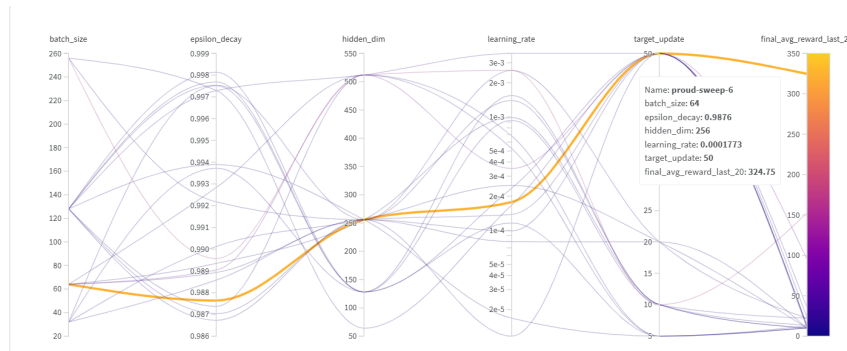


Figure 6: DQN Duelling - Cartpole-v1 (Maximum): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Batch Size:** 64
- **Epsilon decay:** 0.9876
- **Hidden Dimension:** 256
- **Learning rate:** 0.0001773
- **Target Update:** 50

6.3 DQN Duelling - Acrobot-v1: Average

The hyperparameters used for DQN Duelling on Acrobot-v1 (Average) are:

- **Batch Size**
- **Epsilon decay**
- **Hidden Dimension**
- **Learning rate**
- **Target Update**

Figure 7 shows the parallel coordinates plot for the average performance.

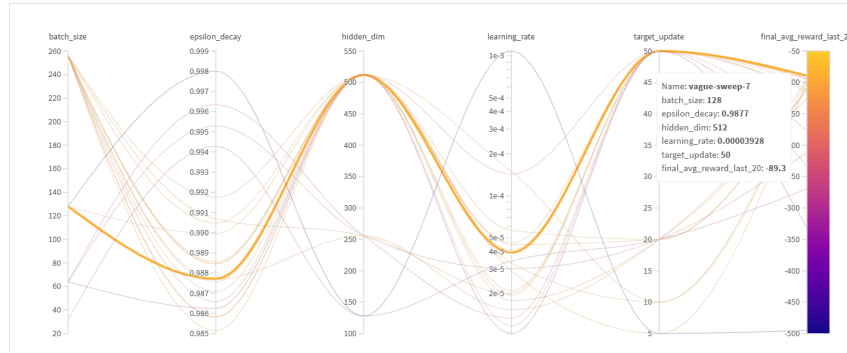


Figure 7: DQN Duelling - Acrobot-v1 (Average): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Batch Size:** 128
- **Epsilon decay:** 0.9877
- **Hidden Dimension:** 512
- **Learning rate:** 0.00003928
- **Target Update:** 50

6.4 DQN Duelling - Acrobot-v1: Maximum

The hyperparameters used for DQN Duelling on Acrobot-v1 (Maximum) are:

- **Batch Size**
- **Epsilon decay**

- **Hidden Dimension**
- **Learning rate**
- **Target Update**

Figure 8 shows the parallel coordinates plot for the maximum performance.

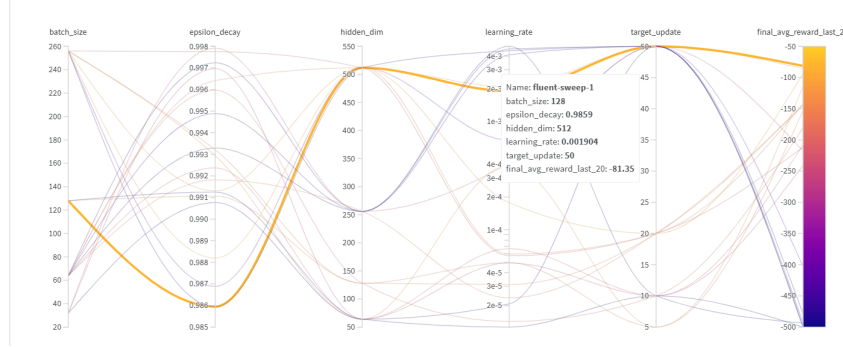


Figure 8: DQN Duelling - Acrobot-v1 (Maximum): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Batch Size:** 128
- **Epsilon decay:** 0.9859
- **Hidden Dimension:** 512
- **Learning rate:** 0.001904
- **Target Update:** 50

6.5 REINFORCE - Cartpole-v1 (No Baseline)

The hyperparameters used for REINFORCE on Cartpole-v1 without a baseline are:

- **Hidden Dimension**
- **Policy Network Learning Rate**

Figure 9 shows the parallel coordinates plot for the no-baseline configuration.

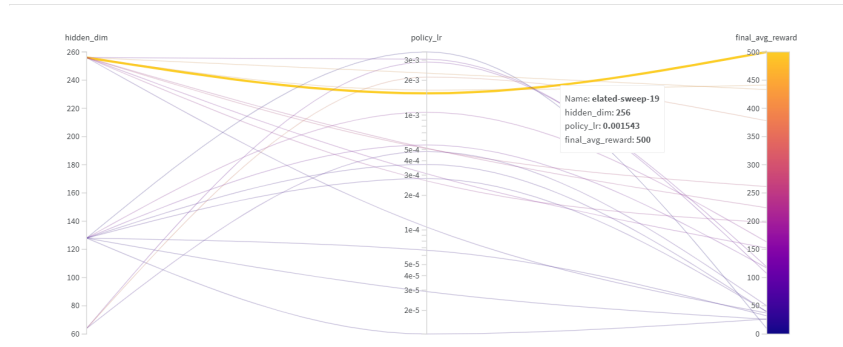


Figure 9: REINFORCE - Cartpole-v1 (No Baseline): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Hidden Dimension:** 256
- **Policy Learning Rate:** 0.001543

6.6 REINFORCE - Cartpole-v1 (Baseline)

The hyperparameters used for REINFORCE on Cartpole-v1 with a baseline are:

- **Hidden Dimension**
- **Policy Learning Rate**
- **Value Network Learning Rate**

Figure 10 shows the parallel coordinates plot for the baseline configuration.

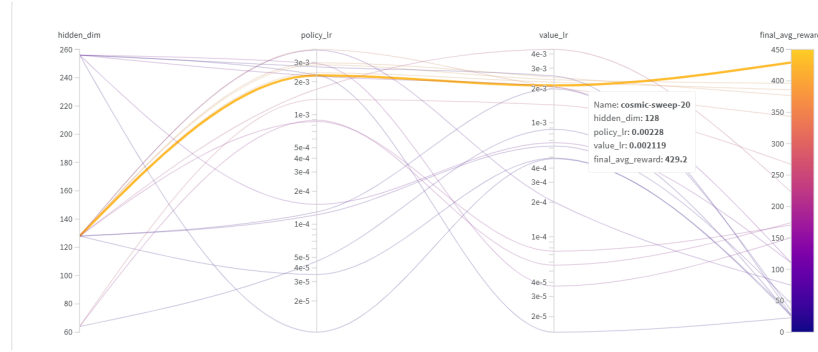


Figure 10: REINFORCE - Cartpole-v1 (Baseline): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Hidden Dimension: 128**
- **Policy Learning Rate: 0.00228**
- **Value Network Learning Rate: 0.002119**

6.7 REINFORCE - Acrobot-v1 (No Baseline)

The hyperparameters used for REINFORCE on Acrobot-v1 without a baseline are:

- **Hidden Dimension**
- **Policy Network Learning Rate**

Figure 11 shows the parallel coordinates plot for the no-baseline configuration.

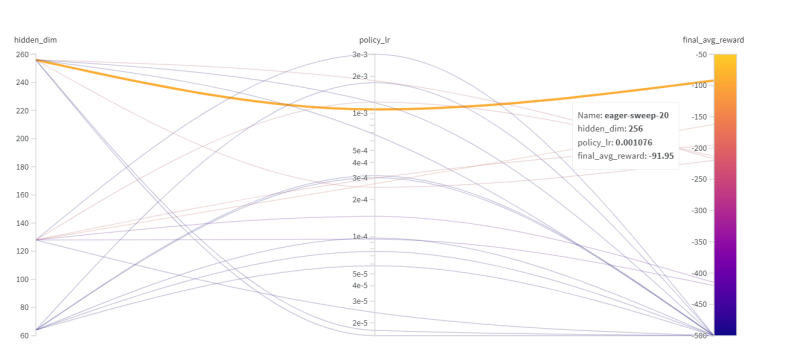


Figure 11: REINFORCE - Acrobot-v1 (No Baseline): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Hidden Dimension: 256**
- **Policy Learning Rate: 0.001076**

6.8 REINFORCE - Acrobot-v1 (Baseline)

The hyperparameters used for REINFORCE on Acrobot-v1 with a baseline are:

- **Hidden Dimension**
- **Policy Learning Rate**
- **Value Network Learning Rate**

Figure 12 shows the parallel coordinates plot for the baseline configuration.

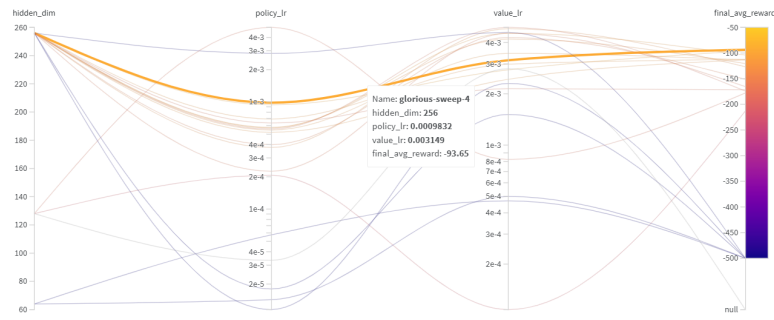


Figure 12: REINFORCE - Acrobot-v1 (Baseline): Parallel coordinates plot.

The best hyperparameter configuration for this experiment was:

- **Hidden Dimension: 256**
- **Policy Learning Rate: 0.0009832**
- **Value Network Learning Rate: 0.003149**

7 Conclusions and Discussions

The experimental results indicate that:

- For Dueling-DQN, the aggregation method (mean vs. max) significantly affects the learning dynamics.
- Max aggregation tends to converge much earlier than average aggregation. But, average aggregation converges as well in 500 episodes for acrobot-v1.
- In Monte-Carlo REINFORCE, incorporating a baseline updated via TD(0) reduces the gradient variance and leads to more stable learning. It can be observed from the graphs how using baseline is much better especially for acrobot-v1.
- Performance trends are consistent across both CartPole-v1 and Acrobot-v1.

8 Code Repository

The complete source code is available on GitHub at: https://github.com/DA6400-RL-JanMay2025/programming-assignment-02-ch21b053_me21b171.