# ASSIGNMENT 4: APPENDIX

**Ozan Ilhan**
s3559513

**Qiu van Leeuwen**
s3547566

## 1 ADVANTAGE ACTOR-CRITIC (A2C)

```python
def A2C_blackjack(model_info, hidden_size, learning_rate, n_repetitions,
    gamma):
    environment = gym.make("Blackjack-v1")
    total_rewards = np.zeros(n_repetitions)

    input_size = len(environment.observation_space)
    output_size = environment.action_space.n

    model = ActorCriticDiscrete(input_size, hidden_size, output_size)
    optimizer_actor = torch.optim.Adam(model.actor.parameters(), lr=
    learning_rate)
    optimizer_critic = torch.optim.Adam(model.critic.parameters(), lr=
    learning_rate)

    wins = np.zeros(n_repetitions)
    draws = np.zeros(n_repetitions)
    losses = np.zeros(n_repetitions)
    for repetition in range(n_repetitions):
        state, _ = environment.reset()
        state = torch.tensor(state, dtype=torch.float32)
        total_reward = 0
        done = False
        count = 0
        while not done:
            count += 1
            probabilities, values = model(state)
            distribution = torch.distributions.Categorical(probabilities)

            action = distribution.sample()

            next_state, reward, done, truncation, _ = environment.step(
    action.item())
            total_rewards[repetition] = reward
            if truncation:
                done = True
            if reward > 0:
                wins[repetition] += 1
            elif reward < 0:
                losses[repetition] += 1
            else:
                draws[repetition] += 1
            next_state = torch.tensor(next_state, dtype=torch.float32)
            _, next_values = model(next_state)
            advantage = reward + (gamma * next_values.detach() * (1 - int
    (done))) - values
            actor_loss = -distribution.log_prob(action) * advantage.
    detach()
            critic_loss = advantage.square()
            total_reward += reward

            optimizer_actor.zero_grad()
            optimizer_critic.zero_grad()

```

```
48            actor_loss.backward()
49            critic_loss.backward()
50
51            optimizer_actor.step()
52            optimizer_critic.step()
53            state = next_state
54    environment.close()
55    torch.save(model.state_dict(), f'A2C_blackjack_{model_info}.pth')
56    return total_rewards
```

Listing 1: The code for the A2C algorithm.

## 2  DEEP Q-LEARNING

```
1  def DQL_blackjack(model_info, hidden_size, learning_rate, n_repetitions,
       gamma):
2      environment = gym.make("Blackjack-v1")
3      total_rewards = np.zeros(n_repetitions)
4
5      input_size = len(environment.observation_space)
6      output_size = environment.action_space.n
7
8      model = DeepQLearning(input_size, hidden_size, output_size)
9      optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
10
11     batch_size = 64
12     memory = deque(maxlen=10000)
13     epsilon = 0.1
14     for repetition in range(n_repetitions):
15         state, _ = environment.reset()
16         state = torch.tensor(state, dtype=torch.float32)
17         total_reward = 0
18         done = False
19         while not done:
20             if np.random.rand() < epsilon:
21                 action = environment.action_space.sample()
22             else:
23                 q_values = model(state)
24                 action = torch.argmax(q_values).item()
25
26             next_state, reward, done, truncation, _ = environment.step(
       action)
27             total_rewards[repetition] = reward
28             if truncation:
29                 done = True
30             next_state = torch.tensor(next_state, dtype=torch.float32)
31             memory.append((state, action, next_state, reward, done))
32             state = next_state
33             total_reward += reward
34
35             if len(memory) > batch_size:
36                 experiences = random.sample(memory, batch_size)
37                 states, actions, next_states, rewards, dones = zip(*
       experiences)
38                 states = np.array(states)
39                 actions = np.array(actions)
40                 rewards = np.array(rewards)
41                 next_states = np.array(next_states)
42                 dones = np.array(dones)
43                 q_values = model(torch.tensor(states, dtype=torch.float32
       ))
44                 next_q_values = model(torch.tensor(next_states, dtype=
       torch.float32))
45
```

```
46                    target_q_values = q_values.clone()
47                    for i in range(len(experiences)):
48                        if dones[i]:
49                            target_q_values[i, actions[i]] = rewards[i]
50                        else:
51                            target_q_values[i, actions[i]] = rewards[i] +
     gamma * torch.max(next_q_values[i]).item()
52                    mse_loss = nn.MSELoss()
53                    loss = mse_loss(q_values, target_q_values.clone().detach
     ().to(dtype=torch.float32))
54                    optimizer.zero_grad()
55                    loss.backward()
56                    optimizer.step()
57                    state = next_state
58          environment.close()
59          torch.save(model.state_dict(), f'DQL_blackjack_{model_info}.pth')
60          return total_rewards
```

Listing 2: The code for the DQL algorithm.