# WIDS Project - Monte Carlo Assignment Week 4: Prediction & Control Solutions

## Monte Carlo Methods for Blackjack

## Introduction

This document provides complete solutions to the Monte Carlo Prediction and Control assignment. We implement algorithms to learn optimal Blackjack strategies through experience, starting from zero knowledge.

# 1 Part I: Environment Setup

## 1.1 Using OpenAI Gymnasium

For this solution, we'll use the Gymnasium library as it provides a standardized Blackjack environment.

```python
import gymnasium as gym
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from collections import defaultdict
from tqdm import tqdm

# Create the Blackjack environment
env = gym.make('Blackjack-v1', sab=True)

# State format: (player_sum, dealer_card, usable_ace)
# Example: (18, 9, False) means player has 18, dealer shows 9, no usable
    ace
```

Listing 1: Environment Setup

# 2 Part II: Algorithm Implementation

## 2.1 Task 1: Monte Carlo Prediction

**Goal:** Estimate the value function $V(s)$ for a fixed policy using First-Visit Monte Carlo.
    **Policy:** Stick if sum is 20 or 21, otherwise hit.

```python
def simple_policy(state):
    """
    Simple policy: Stick on 20 or 21, otherwise hit.

    Args:
        state: Tuple of (player_sum, dealer_card, usable_ace)

    Returns:
        action: 0 (stick) or 1 (hit)
    """
    player_sum = state[0]
    if player_sum >= 20:
        return 0  # Stick
    else:
        return 1  # Hit


def monte_carlo_prediction(env, policy, num_episodes=10000):
    """
    First-Visit Monte Carlo Prediction.

    Estimates V(s) by averaging returns from first visits to each state.

    Args:
        env: Gymnasium environment
        policy: Function that takes state and returns action
        num_episodes: Number of episodes to run

    Returns:
        V: Dictionary mapping states to their estimated values
    """
    # Store sum of returns and visit counts for each state
    returns_sum = defaultdict(float)
    returns_count = defaultdict(int)
    V = defaultdict(float)

    for episode in tqdm(range(num_episodes), desc="MC Prediction"):
        # Generate an episode
        episode_data = []
        state, info = env.reset()
        done = False

        while not done:
            action = policy(state)
            next_state, reward, terminated, truncated, info = env.step(
    action)
            episode_data.append((state, reward))
            state = next_state
            done = terminated or truncated

        # Track which states we've already seen in this episode
        visited_states = set()

        # Calculate returns for each state (working backwards)
```

```
54          G = 0   # Return
55          for state, reward in reversed(episode_data):
56              G = reward   # In Blackjack, only terminal reward matters
57
58              # First-visit: only update if we haven't seen this state
    before
59              if state not in visited_states:
60                  visited_states.add(state)
61                  returns_sum[state] += G
62                  returns_count[state] += 1
63                  V[state] = returns_sum[state] / returns_count[state]
64
65      return V
66
67
68 # Run Monte Carlo Prediction
69 print("Running Monte Carlo Prediction...")
70 V = monte_carlo_prediction(env, simple_policy, num_episodes=10000)
71
72 # Print results for specific states
73 print("\n=== Value Function Results ===")
74 # Find states with player sum 21 and 5
75 states_21 = [s for s in V.keys() if s[0] == 21]
76 states_5 = [s for s in V.keys() if s[0] == 5]
77
78 if states_21:
79     avg_value_21 = np.mean([V[s] for s in states_21])
80     print(f"Average Value of Player Sum = 21: {avg_value_21:.4f}")
81     print(f"  (This is high because 21 usually wins!)")
82 else:
83     print("No states with player sum 21 visited")
84
85 if states_5:
86     avg_value_5 = np.mean([V[s] for s in states_5])
87     print(f"Average Value of Player Sum = 5: {avg_value_5:.4f}")
88     print(f"  (This is low because 5 requires many hits, risking bust)")
89 else:
90     print("No states with player sum 5 visited")
```

Listing 2: Monte Carlo Prediction Implementation

**Explanation:**

- We run 10,000 episodes following the simple policy

- For each state visited, we track the return (reward) received

- Using First-Visit MC, we only count the first time a state appears in an episode

- $V(s) =$ average of all returns from that state

- Player sum of 21 has high value (usually wins)

- Player sum of 5 has low value (far from 21, many hits needed)

3

## 2.2 Task 2: Monte Carlo Control

**Goal:** Find the optimal policy $\pi^*$ by learning action-values $Q(s, a)$.

```python
def epsilon_greedy_policy(Q, state, epsilon, n_actions=2):
    """
    Epsilon-greedy policy: explore with probability epsilon,
    exploit (choose best action) with probability 1-epsilon.

    Args:
        Q: Action-value function (dictionary)
        state: Current state
        epsilon: Exploration probability
        n_actions: Number of possible actions

    Returns:
        action: Selected action (0 or 1)
    """
    if np.random.random() < epsilon:
        # Explore: choose random action
        return np.random.randint(n_actions)
    else:
        # Exploit: choose best action
        q_values = [Q[(state, a)] for a in range(n_actions)]
        return np.argmax(q_values)


def monte_carlo_control(env, num_episodes=500000, alpha=0.01,
                        epsilon_start=1.0, epsilon_end=0.1, epsilon_decay
    =0.99999):
    """
    Monte Carlo Control with epsilon-greedy exploration.

    Learns optimal policy by updating Q(s,a) values based on returns.

    Args:
        env: Gymnasium environment
        num_episodes: Number of training episodes
        alpha: Learning rate for Q-value updates
        epsilon_start: Initial exploration rate
        epsilon_end: Final exploration rate
        epsilon_decay: Rate of epsilon decay

    Returns:
        Q: Learned action-value function
        rewards_history: List of rewards per episode
        policy: Optimal policy extracted from Q
    """
    # Initialize Q(s,a) to zero for all state-action pairs
    Q = defaultdict(float)

    # Track rewards for plotting learning curve
    rewards_history = []
    epsilon = epsilon_start
```

```python
    for episode in tqdm(range(num_episodes), desc="MC Control"):
        # Generate an episode using epsilon-greedy policy
        episode_data = []
        state, info = env.reset()
        done = False
        episode_reward = 0

        while not done:
            action = epsilon_greedy_policy(Q, state, epsilon)
            next_state, reward, terminated, truncated, info = env.step(
    action)
            episode_data.append((state, action, reward))
            episode_reward += reward
            state = next_state
            done = terminated or truncated

        rewards_history.append(episode_reward)

        # Update Q-values using returns from this episode
        G = 0  # Return
        visited_state_actions = set()

        for state, action, reward in reversed(episode_data):
            G = reward  # In Blackjack, only terminal reward

            # First-visit update
            if (state, action) not in visited_state_actions:
                visited_state_actions.add((state, action))

                # Incremental update formula
                Q[(state, action)] = Q[(state, action)] + alpha * (G - Q[(
    state, action)])

        # Decay epsilon (explore less over time)
        epsilon = max(epsilon_end, epsilon * epsilon_decay)

    # Extract optimal policy from Q-values
    policy = {}
    states_visited = set(s for (s, a) in Q.keys())
    for state in states_visited:
        q_values = [Q[(state, a)] for a in range(2)]
        policy[state] = np.argmax(q_values)

    return Q, rewards_history, policy


# Run Monte Carlo Control
print("\nRunning Monte Carlo Control...")
Q, rewards_history, optimal_policy = monte_carlo_control(
    env,
    num_episodes=500000,
    alpha=0.01,
    epsilon_start=1.0,
    epsilon_end=0.1,
```

```
103      epsilon_decay =0.99999
104  )
105
106  print (f"\nLearned Q-values for {len(Q)} state-action pairs")
107  print (f"Derived optimal policy for {len(optimal_policy)} states")
```
Listing 3: Monte Carlo Control Implementation

**Explanation:**

- **Initialization:** Start with $Q(s, a) = 0$ for all states and actions

- **Epsilon-greedy:** Balance exploration (random actions) and exploitation (best known actions)

- **Episode generation:** Play through game using current policy

- **Q-value update:** $Q(s, a) \leftarrow Q(s, a) + \alpha[G_t - Q(s, a)]$

- **Epsilon decay:** Gradually reduce exploration as we learn

- Over 500,000 episodes, the agent learns which actions lead to wins

## 2.3  Task 3: Visualization

```
1  def plot_learning_curve (rewards_history , window =1000):
2      """
3      Plot the rolling average reward to show learning progress.
4
5      Args:
6          rewards_history: List of rewards from each episode
7          window: Size of rolling average window
8      """
9      # Calculate rolling average
10     rolling_avg = np.convolve (rewards_history ,
11                               np.ones (window)/window ,
12                               mode ='valid')
13
14     plt.figure (figsize =(12, 6))
15     plt.plot (rolling_avg , linewidth =2)
16     plt.xlabel ('Episode', fontsize =12)
17     plt.ylabel (f'Rolling Average Reward (window={window})', fontsize =12)
18     plt.title ('Monte Carlo Control: Learning Curve', fontsize =14,
     fontweight ='bold')
19     plt.grid (True , alpha =0.3)
20     plt.axhline (y=0, color ='r', linestyle ='--', alpha =0.5, label ='Break-
     even')
21     plt.legend ()
22     plt.tight_layout ()
23     plt.savefig ('learning_curve.png', dpi =300, bbox_inches ='tight')
24     plt.show ()
25
```

```python
26      print(f"Final average reward: {np.mean(rewards_history[-10000:]):.4f}"
    )


29  def plot_strategy_heatmap(policy):
30      """
31      Create a heatmap showing the optimal action for each state.

33      Args:
34          policy: Dictionary mapping states to optimal actions
35      """
36      # Create separate heatmaps for with/without usable ace
37      fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))

39      # Prepare data grids
40      player_range = range(12, 22)  # Player sum from 12 to 21
41      dealer_range = range(1, 11)   # Dealer card from 1 (Ace) to 10

43      for usable_ace, ax, title in [(False, ax1, 'No Usable Ace'),
44                                    (True, ax2, 'Usable Ace')]:
45          strategy_grid = np.zeros((len(player_range), len(dealer_range)))

47          for i, player_sum in enumerate(player_range):
48              for j, dealer_card in enumerate(dealer_range):
49                  state = (player_sum, dealer_card, usable_ace)
50                  if state in policy:
51                      strategy_grid[i, j] = policy[state]
52                  else:
53                      strategy_grid[i, j] = np.nan

55          # Create heatmap
56          sns.heatmap(strategy_grid,
57                      ax=ax,
58                      cmap=['darkred', 'darkgreen'],
59                      cbar_kws={'label': 'Action',
60                                'ticks': [0.25, 0.75]},
61                      xticklabels=['A'] + list(range(2, 11)),
62                      yticklabels=player_range,
63                      linewidths=0.5,
64                      linecolor='gray',
65                      vmin=0, vmax=1)

67          ax.set_xlabel('Dealer Showing', fontsize=12, fontweight='bold')
68          ax.set_ylabel('Player Sum', fontsize=12, fontweight='bold')
69          ax.set_title(title, fontsize=14, fontweight='bold')

71          # Fix colorbar labels
72          cbar = ax.collections[0].colorbar
73          cbar.set_ticklabels(['STICK (0)', 'HIT (1)'])

75      plt.suptitle('Optimal Blackjack Strategy',
76                   fontsize=16, fontweight='bold', y=1.02)
77      plt.tight_layout()
78      plt.savefig('strategy_heatmap.png', dpi=300, bbox_inches='tight')
```

```
79      plt.show()
80
81
82  # Generate visualizations
83  print("\nGenerating visualizations...")
84  plot_learning_curve(rewards_history, window=1000)
85  plot_strategy_heatmap(optimal_policy)
```

Listing 4: Visualization Code

**Interpretation:**

- **Learning Curve:** Shows improvement over time. Should trend upward and stabilize.

- **Strategy Heatmap:**

    - Green = HIT (action 1)
    - Red = STICK (action 0)
    - Generally: stick on high sums (17+), hit on low sums
    - Strategy differs with/without usable ace

# 3 Part III: Conceptual Questions

## 3.1 Question A: The Infinite Deck Assumption

### 3.1.1 1. Non-Stationarity and Markov Property

**Answer:**

In real casinos with finite decks (6-8 deck shoe), the environment becomes **non-Markovian** if we only use (PlayerSum, DealerCard) as the state.

**Why?**

The **Markov Property** states that the future depends only on the current state, not on the history. Mathematically:

$$P(S_{t+1}|S_t, S_{t-1}, \ldots, S_0) = P(S_{t+1}|S_t) \tag{1}$$

With a finite deck:

- Cards already dealt affect future probabilities

- If many 10-value cards have been dealt, the probability of drawing another 10 decreases

- The state (18, 10) has different implications early vs. late in the shoe

- **Example:** If 30 face cards have been dealt, hitting on 12 is safer than if only 5 face cards have been dealt

Therefore, (PlayerSum, DealerCard) alone does **not** contain all information needed to predict the future, violating the Markov property.

8

### 3.1.2  2. State Space Design for Card Counting

**Answer:**

To enable card counting, we need to augment the state to capture information about remaining cards:

**Option 1: Running Count**

$$S = (\text{PlayerSum}, \text{DealerCard}, \text{UsableAce}, \text{RunningCount}) \tag{2}$$

Where RunningCount is:

- +1 for each low card dealt (2-6)

- 0 for each neutral card (7-9)

- -1 for each high card dealt (10, J, Q, K, A)

**Option 2: True Count**

$$S = (\text{PlayerSum}, \text{DealerCard}, \text{UsableAce}, \text{TrueCount}, \text{DecksRemaining}) \tag{3}$$

Where: $\text{TrueCount} = \frac{\text{RunningCount}}{\text{DecksRemaining}}$

**Option 3: Full Deck Composition (Most Markovian)**

$$S = (\text{PlayerSum}, \text{DealerCard}, \text{UsableAce}, \mathbf{c}) \tag{4}$$

Where $\mathbf{c} = (c_1, c_2, \ldots, c_{10})$ is a vector of remaining card counts.

### 3.1.3  3. Trade-off: Convergence Speed

**Answer:**

Expanding the state space creates a fundamental trade-off:

| Pros | Cons |
| --- | --- |
| More Markovian (better) | Larger state space |
| Better optimal policy | Slower convergence |
| Can exploit card counting | Need more episodes |
| Higher theoretical reward | More memory required |

**Mathematical Explanation:**

With the original state space of size $|S_1| \approx 200$ (combinations of player sum, dealer card, usable ace), adding a running count from -50 to +50 increases the state space to:

$$|S_2| = |S_1| \times 100 \approx 20,000 \text{ states} \tag{5}$$

Since Monte Carlo methods require visiting each state-action pair multiple times:

- Episodes needed grows roughly with $|S|$

- 100× more states $\Rightarrow$ 100× more episodes needed

- Each state visited less frequently

- Slower learning, longer training time

**Practical Solution:** Use function approximation (neural networks) instead of tabular methods for large state spaces.

## 3.2 Question B: First-Visit vs. Every-Visit

**Scenario:** $A \rightarrow B \rightarrow A \rightarrow$ Terminate
**Rewards:**

- $R(A \rightarrow B) = +1$

- $R(B \rightarrow A) = -1$

- $R(A \rightarrow \text{Term}) = +10$

### 3.2.1 1. First-Visit Monte Carlo

**Solution:**
In First-Visit MC, we only count the **first time** state $A$ appears.
**Calculation:**

$$\begin{align}
G_{\text{first}} &= R_1 + R_2 + R_3 \tag{6} \\
&= (+1) + (-1) + (+10) \tag{7} \\
&= +10 \tag{8}
\end{align}$$

The return from the first visit to $A$ is $+\mathbf{10}$.
After many episodes, $V(A) \approx$ average of all such first-visit returns.

### 3.2.2 2. Every-Visit Monte Carlo

**Solution:**
In Every-Visit MC, we count **both** times state $A$ appears.
**Visit 1 (first occurrence):**

$$G_1 = (+1) + (-1) + (+10) = +10 \tag{9}$$

**Visit 2 (second occurrence):**

$$G_2 = (+10) = +10 \tag{10}$$

For this episode, state $A$ contributes two returns: $+10$ and $+10$.
After many episodes, $V(A) \approx$ average of **all** returns from visiting $A$.

### 3.2.3 3. Comparison

| Method | Returns from A | Average |
|---|---|---|
| First-Visit MC | $\{+10\}$ | $+10$ |
| Every-Visit MC | $\{+10, +10\}$ | $+10$ |

In this particular example, both methods give the same result. However, in general:

- Every-Visit uses more data (all visits)

- First-Visit has theoretical guarantees under weaker conditions

- Both converge to the true value function $V(s)$ as episodes $\rightarrow \infty$

# 4 Complete Code Solution

Here is the complete, runnable code combining all components:

```python
import gymnasium as gym
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from collections import defaultdict
from tqdm import tqdm

# ================================================
# PART 1: MONTE CARLO PREDICTION
# ================================================

def simple_policy(state):
    """Stick on 20-21, otherwise hit."""
    return 0 if state[0] >= 20 else 1

def monte_carlo_prediction(env, policy, num_episodes=10000):
    """First-Visit MC Prediction."""
    returns_sum = defaultdict(float)
    returns_count = defaultdict(int)
    V = defaultdict(float)

    for _ in tqdm(range(num_episodes), desc="MC Prediction"):
        episode_data = []
        state, _ = env.reset()
        done = False

        while not done:
            action = policy(state)
            next_state, reward, terminated, truncated, _ = env.step(action
    )
            episode_data.append((state, reward))
            state = next_state
            done = terminated or truncated

        visited_states = set()
        G = 0

        for state, reward in reversed(episode_data):
            G = reward
            if state not in visited_states:
                visited_states.add(state)
                returns_sum[state] += G
                returns_count[state] += 1
                V[state] = returns_sum[state] / returns_count[state]

    return V

# ================================================
# PART 2: MONTE CARLO CONTROL
# ================================================
```

```python
50
51  def epsilon_greedy_policy(Q, state, epsilon, n_actions=2):
52      """Epsilon-greedy action selection."""
53      if np.random.random() < epsilon:
54          return np.random.randint(n_actions)
55      else:
56          q_values = [Q[(state, a)] for a in range(n_actions)]
57          return np.argmax(q_values)
58
59  def monte_carlo_control(env, num_episodes=500000, alpha=0.01):
60      """MC Control with epsilon-greedy."""
61      Q = defaultdict(float)
62      rewards_history = []
63      epsilon = 1.0
64      epsilon_min = 0.1
65      epsilon_decay = 0.99999
66
67      for _ in tqdm(range(num_episodes), desc="MC Control"):
68          episode_data = []
69          state, _ = env.reset()
70          done = False
71          episode_reward = 0
72
73          while not done:
74              action = epsilon_greedy_policy(Q, state, epsilon)
75              next_state, reward, terminated, truncated, _ = env.step(action
    )
76              episode_data.append((state, action, reward))
77              episode_reward += reward
78              state = next_state
79              done = terminated or truncated
80
81          rewards_history.append(episode_reward)
82
83          visited = set()
84          G = 0
85
86          for state, action, reward in reversed(episode_data):
87              G = reward
88              if (state, action) not in visited:
89                  visited.add((state, action))
90                  Q[(state, action)] += alpha * (G - Q[(state, action)])
91
92          epsilon = max(epsilon_min, epsilon * epsilon_decay)
93
94      # Extract policy
95      policy = {}
96      for (state, action) in Q.keys():
97          if state not in policy:
98              q_vals = [Q[(state, a)] for a in range(2)]
99              policy[state] = np.argmax(q_vals)
100
101     return Q, rewards_history, policy
102
```

```
103  # ==============================================
104  # MAIN EXECUTION
105  # ==============================================
106
107  if __name__ == "__main__":
108      # Setup
109      env = gym.make('Blackjack-v1', sab=True)
110
111      # Task 1: Prediction
112      print("Task 1: Monte Carlo Prediction")
113      V = monte_carlo_prediction(env, simple_policy, 10000)
114
115      states_21 = [s for s in V.keys() if s[0] == 21]
116      states_5 = [s for s in V.keys() if s[0] == 5]
117
118      if states_21:
119          print(f"V(player_sum=21): {np.mean([V[s] for s in states_21]):.4f}
     ")
120      if states_5:
121          print(f"V(player_sum=5): {np.mean([V[s] for s in states_5]):.4f}")
122
123      # Task 2: Control
124      print("\nTask 2: Monte Carlo Control")
125      Q, rewards, policy = monte_carlo_control(env, 500000)
126      print(f"Learned policy for {len(policy)} states")
127      print(f"Final avg reward: {np.mean(rewards[-10000:]):.4f}")
128
129      env.close()
```

Listing 5: Complete Solution Script

# 5 Conclusion

This assignment demonstrates the power of Monte Carlo methods:

- **Prediction:** Evaluates a fixed policy by averaging returns

- **Control:** Learns optimal policy through exploration and exploitation

- **Model-free:** No knowledge of transition probabilities needed

- **Experience-based:** Learns purely from playing the game

The agent starts with zero knowledge and discovers the optimal Blackjack strategy through trial and error!