

-----Mandatory Information to fill-----

Group ID: 88

Group Members Name with Student ID:

1. Student 1: SIGINAM SIVASAI, ID: 2023AA05371
2. Student 2: PEYALA SAMARASIMHA REDDY, ID: 2023AA05072
3. Student 3: PEGALLAPATI SAI MAHARSHI, ID: 2023AA05924
4. Student 4: CHADALAWADA VISWANATH HEMANTH, ID: 2023AA05195

-----Write your remarks (if any) that you want should get
consider at the time of evaluation-----

Remarks: ##Add here

Problem Statement

Develop a reinforcement learning agent using dynamic programming methods to solve the Dice game optimally. The agent will learn the optimal policy by iteratively evaluating and improving its strategy based on the state-value function and the Bellman equations.

Scenario:

A player rolls a 6-sided die with the objective of reaching a score of **exactly** 100. On each turn, the player can choose to stop and keep their current score or continue rolling the die. If the player rolls a 1, they lose all points accumulated in that turn and the turn ends. If the player rolls any other number (2-6), that number is added to their score for that turn. The game ends when the player decides to stop and keep their score OR when the player's score reaches 100. The player wins if they reach a score of exactly 100, and loses if they roll a 1 when their score is below 100.

Environment Details

- The environment consists of a player who can choose to either roll a 6-sided die or stop at any point.
- The player starts with an initial score (e.g., 0) and aims to reach a score of exactly 100.
- If the player rolls a 1, they lose all points accumulated in that turn and the turn ends. If they roll any other number (2-6), that number is added to their score for that turn.
- The goal is to accumulate a total of exactly 100 points to win, or to stop the game before reaching 100 points.

States

- State s : Represents the current score of the player, ranging from 0 to 100.
- Terminal States:
 - State $s = 100$: Represents the player winning the game by reaching the goal of 100 points.

- State $s = 0$: Represents the player losing all points accumulated in the turn due to rolling a 1.

Actions

- Action a: Represents the decision to either "roll" the die or "stop" the game at the current score.
- The possible actions in any state s are either "roll" or "stop".

Expected Outcomes:

1. Use dynamic programming methods value iteration, policy improvement and policy evaluation to find the optimal policy for the Dice Game.
2. Implement an epsilon-greedy policy for action selection during training to balance exploration and exploitation.
3. Evaluate the agent's performance in terms of the probability of reaching exactly 100 points after learning the optimal policy.
4. Use the agent's policy as the best strategy for different betting scenarios within the problem.

Code Execution

Initialize constants

```
In [3]: # Importing the necessary libraries
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [4]: # Constants
goal = 100
gamma = 1.0
prob_roll = 1/6

# Initialize value function and policy
V = np.zeros(goal + 1)
policy = np.zeros(goal + 1, dtype=int) # 0 for "stop", 1 for "roll"
```

Design a DiceGame Environment (1M)

```
In [5]: # Code for Dataset loading and print dataset statistics along with reward function
#-----write your code below this line-----

# Define the environment class for the Dice game
class DiceGameEnvironment:
    def __init__(self, goal=100):
        self.goal = goal
        self.state = 0

    def reset(self):
        self.state = 0
        return self.state

    def step(self, action):
        if action == "stop":
```

```

        return self.state, 0, True

    roll = np.random.randint(1, 7)
    if roll == 1:
        return 0, -self.state, True
    else:
        self.state += roll
        if self.state == self.goal:
            return self.state, self.goal, True
        return self.state, roll, False

```

Define reward function

```

In [6]: #Calculate reward function for 'stop' and 'roll' actions
#-----write your code below this line-----
def reward_function(state, action):
    if action == "stop":
        return state if state == goal else -state
    else:
        return 0

```

Policy Iteration Function Definition (0.5M)

```

In [7]: #For each state, Store old_policy of state s.
#Determine best_action based on maximum reward. Update policy[s] to best_action.
#Return stable when old policy = policy[s]

#-----write your code below this line-----
def policy_iteration(env, gamma=1.0, theta=1e-6):
    policy = np.zeros(goal + 1, dtype=int)
    V = np.zeros(goal + 1)
    stable = False

    while not stable:
        # Policy Evaluation
        while True:
            delta = 0
            for s in range(goal + 1):
                v = V[s]
                if policy[s] == 0: # stop
                    V[s] = reward_function(s, "stop")
                else: # roll
                    V[s] = sum(prob_roll * (reward_function(s, "roll") + gamma * V[min(s + i, goal)]
                                for i in range(2, 7)))
            delta = max(delta, abs(v - V[s]))
            if delta < theta:
                break

        # Policy Improvement
        stable = True
        for s in range(goal + 1):
            old_action = policy[s]
            action_values = [
                reward_function(s, "stop"),
                sum(prob_roll * (reward_function(s, "roll") + gamma * V[min(s + i, goal)]
                                for i in range(2, 7)))
            ]
            policy[s] = np.argmax(action_values)
            if old_action != policy[s]:
                stable = False
    return policy, V

```

Value Iteration Function Definition (0.5M)

```
In [8]: # Iterate over all states except terminal state untill convergence

# Calculate expected returns  $V(s)$  for current policy by considering all possible actions

#If action is stop:
    #Calculate reward for stopping and append to rewards.
#If action is roll:
    #For each possible roll outcome (1 to 6), Determine next_s based on roll.

# Update  $V(s)$  using the Bellman equation.

#Determine max_reward from rewards
#With probability epsilon, randomly choose a reward from rewards.

#Check convergence if delta is less than a small threshold.

#-----write your code below this line-----

def value_iteration(env, gamma=1.0, theta=1e-6):
    V = np.zeros(goal + 1)
    policy = np.zeros(goal + 1, dtype=int)

    while True:
        delta = 0
        for s in range(goal + 1):
            v = V[s]
            action_values = [
                reward_function(s, "stop"),
                sum(prob_roll * (reward_function(s, "roll") + gamma * V[min(s + i, goal)]
                )
            ]
            V[s] = max(action_values)
            delta = max(delta, abs(v - V[s]))
        if delta < theta:
            break

    for s in range(goal + 1):
        action_values = [
            reward_function(s, "stop"),
            sum(prob_roll * (reward_function(s, "roll") + gamma * V[min(s + i, goal)]
            )
        ]
        policy[s] = np.argmax(action_values)

    return policy, V
```

Executing Policy Iteration and Value Iteration Functions (1M)

Print all the iterations for both Policy and Value Iteration approaches separately. (Mandatory)

```
In [12]: #Simulate the game for 100 states. Use the learned policy to get the actions.
#when its roll, randomly generate a number to find the reward.
#when its stop, get the respective reward
#determine the total cumulative reward

#-----write your code below this line-----
# Instantiate the environment
env = DiceGameEnvironment(goal)

# Policy Iteration
```

```
optimal_policy_pi, optimal_value_pi = policy_iteration(env)
print("Policy Iteration Optimal Policy:", optimal_policy_pi)
print("Policy Iteration Optimal Value Function:", optimal_value_pi)

# Value Iteration
optimal_policy_vi, optimal_value_vi = value_iteration(env)
print("Value Iteration Optimal Policy:", optimal_policy_vi)
print("Value Iteration Optimal Value Function:", optimal_value_vi)
```

[illegible]

Print the Learned Optimal Policy, Optimal Value Function (0.5M)

```
In [14]: #-----write your code below this line-----
```

```
Learned Optimal Policy (Policy Iteration): [1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
Learned Optimal Value Function (Policy Iteration): [ 1.01969827    1.06671037    1.1158  
8992     1.16733683     1.22115566  
   1.27745574     1.33635147     1.39796253     1.46241411     1.52983715  
   1.60036867     1.67415196     1.75133695     1.83208048     1.9165466  
   2.00490694     2.09734104     2.19403671     2.29519044     2.40100775  
   2.51170365     2.62750307     2.7486413      2.87536448     3.0079301  
   3.14660752     3.2916785       3.44343784     3.60219387     3.76826918  
   3.94200122     4.12374299     4.31386377     4.51274986     4.72080539  
   4.93845309     5.16613521     5.40431438     5.65347454     5.91412195  
   6.18678624     6.47202144     6.77040712     7.08254953     7.40908289  
   7.75067072     8.10800718     8.4818183       8.87286362     9.28193737  
   9.70987088    10.15753414    10.62583706    11.11573034    11.62820929  
  12.16431336    12.72513523    13.31181665    13.92554785    14.56756893  
  15.23918707    15.94175965    16.67674786    17.44563637    18.24995617  
  19.09131354    19.97146847    20.89218337    21.85556564    22.86328721  
  23.91723235    25.0196127     26.17311292    27.37985502    28.64358083  
  29.96356177    31.34328356    32.78739501    34.30085634    35.88403347  
  37.54591658    39.26316921    41.06572574    42.96552506    44.96480148  
  47.04497933    49.23446788    51.36924154    53.7808642     56.36359739  
  59.04063786    61.71553498    64.50617284    66.58950617    70.83333333  
  74.53703704    77.77777778    80.55555556    83.33333333    83.33333333  
100.]  
Learned Optimal Policy (Value Iteration): [1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
Learned Optimal Value Function (Value Iteration): [ 1.01969823    1.06671036    1.11588  
991     1.16733683     1.22115565  
   1.27745574     1.33635147     1.39796253     1.46241411     1.52983715  
   1.60036867     1.67415196     1.75133695     1.83208048     1.9165466  
   2.00490694     2.09734104     2.19403671     2.29519044     2.40100775  
   2.51170365     2.62750307     2.7486413       2.87536448     3.0079301  
   3.14660752     3.2916785       3.44343784     3.60219387     3.76826918  
   3.94200122     4.12374299     4.31386377     4.51274986     4.72080539  
   4.93845309     5.16613521     5.40431438     5.65347454     5.91412195  
   6.18678624     6.47202144     6.77040712     7.08254953     7.40908289  
   7.75067072     8.10800718     8.4818183       8.87286362     9.28193737  
   9.70987088    10.15753414    10.62583706    11.11573034    11.62820929  
  12.16431336    12.72513523    13.31181665    13.92554785    14.56756893  
  15.23918707    15.94175965    16.67674786    17.44563637    18.24995617  
  19.09131354    19.97146847    20.89218337    21.85556564    22.86328721  
  23.91723235    25.0196127     26.17311292    27.37985502    28.64358083  
  29.96356177    31.34328356    32.78739501    34.30085634    35.88403347  
  37.54591658    39.26316921    41.06572574    42.96552506    44.96480148  
  47.04497933    49.23446788    51.36924154    53.7808642     56.36359739  
  59.04063786    61.71553498    64.50617284    66.58950617    70.83333333  
  74.53703704    77.77777778    80.55555556    83.33333333    83.33333333  
100.]
```

Change in environment details (1M)

Consider the following scenario:

- ```
In [16]: #-----write your code below this line-----
1. What happens if we change the goal score to 50 instead of 100? How does it affect the optimal policy and value function?
env_50 = DiceGameEnvironment(goal=50)
optimal_policy_50, optimal_value_50 = value_iteration(env_50)
print("Value Iteration Optimal Policy with goal 50:", optimal_policy_50)
print("Value Iteration Optimal Value Function with goal 50:", optimal_value_50)
```

```
In [18]: # 2. How would the optimal policy and value function change if the die had 8 sides in
def reward_function_8_sides(state, action):
 if action == "stop":
 return state if state == goal else -state
 else:
 return 0

def value_iteration_8_sides(env, gamma=1.0, theta=1e-6):
 V = np.zeros(goal + 1)
 policy = np.zeros(goal + 1, dtype=int)
 prob_roll = 1/8

 while True:
 delta = 0
 for s in range(goal + 1):
 v = V[s]
 action_values = [
 reward_function_8_sides(s, "stop"),
 sum(prob_roll * (reward_function_8_sides(s, "roll") + gamma * V[min(s
```

```

Value Iteration Optimal Policy with 8-sided die: [1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1]
Value Iteration Optimal Value Function with 8-sided die: [3.50834557 3.62548594
3.7465375 3.87163083 4.0009009
4.13448717 4.27253375 4.41518957 4.56260853 4.71494967
4.87237733 5.03506136 5.20317724 5.37690635 5.55643612
5.74196021 5.93367877 6.13179863 6.33653353 6.54810432
6.76673927 6.99267422 7.22615293 7.46742727 7.71675752
7.97441268 8.2406707 8.51581881 8.80015386 9.09398258
9.39762197 9.71139959 10.03565394 10.37073483 10.71700373
11.07483423 11.44461237 11.82673707 12.22162054 12.62968875
13.05138191 13.48715498 13.93747822 14.40283747 14.88373467
15.38068834 15.89423444 16.42492714 16.97333958 17.54006408
18.12571155 18.73091225 19.35631773 20.00260316 20.67046874
21.36063915 22.07386004 22.81089132 23.57251785 24.35956158
25.17288659 26.01339343 26.88200237 27.77962715 28.70714157
29.66553009 30.65591144 31.6794867 32.73744812 33.8308739
34.96062537 36.12725697 37.33263824 38.57896223 39.86808875
41.20113952 42.57828013 43.99863708 45.46030981 46.97568837
48.54955421 50.18110084 51.86554568 53.59540508 55.36149269
57.15369163 59.09871683 61.14048097 63.23347386 65.34110437
67.43428027 69.49019358 71.49128318 74.65891838 77.47459412
79.97741699 82.20214844 84.1796875 85.9375 87.5
100.]

```

```
In [19]: # 3. Experiment with different discount factors (e.g., 0.9, 0.95). How does discounti
Experiment with different discount factors
gamma_values = [0.9, 0.95]
for gamma in gamma_values:
 optimal_policy_gamma, optimal_value_gamma = value_iteration(env, gamma=gamma)
 print(f"Value Iteration Optimal Policy with gamma {gamma}:", optimal_policy_gamma)
 print(f"Value Iteration Optimal Value Function with gamma {gamma}:", optimal_valu
```





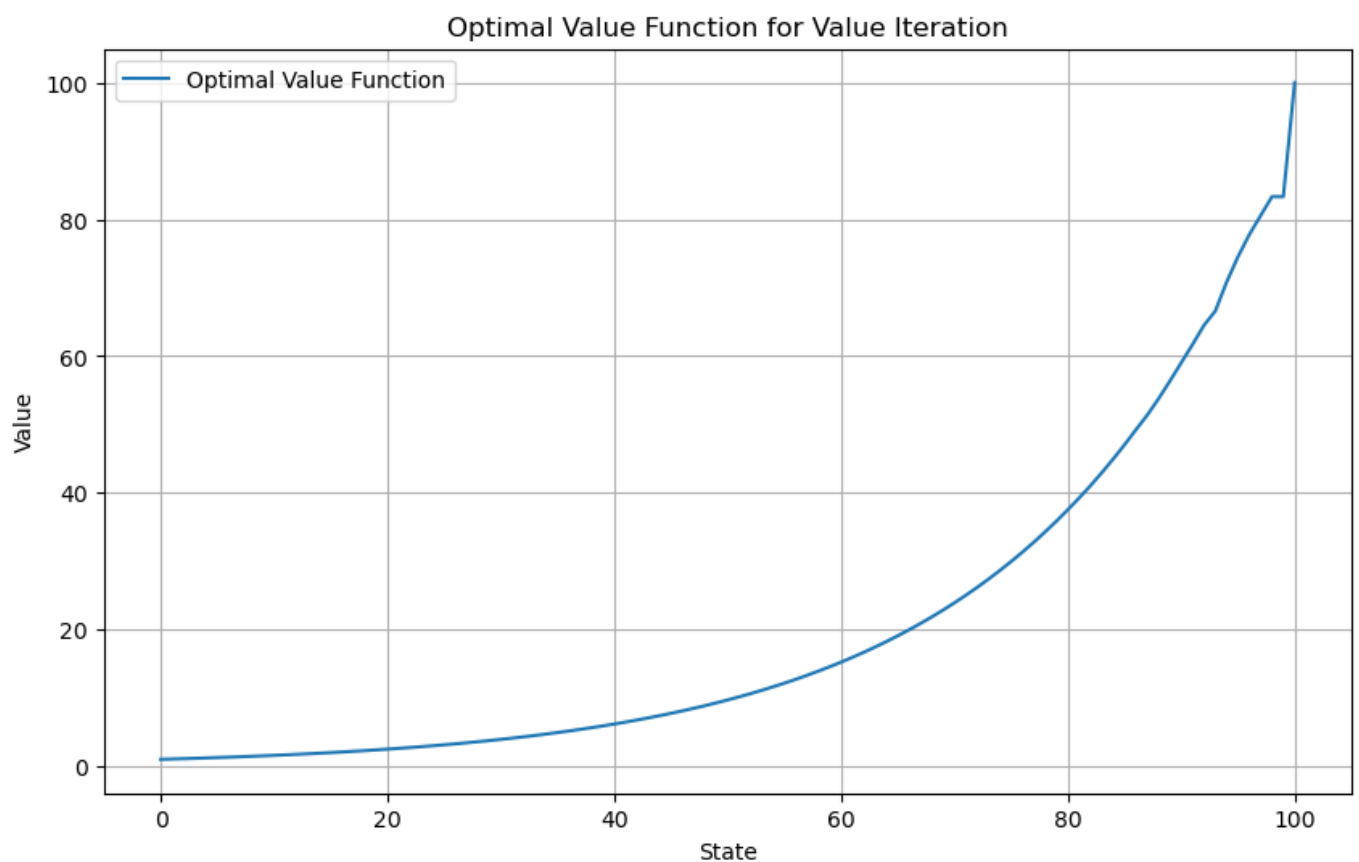
```

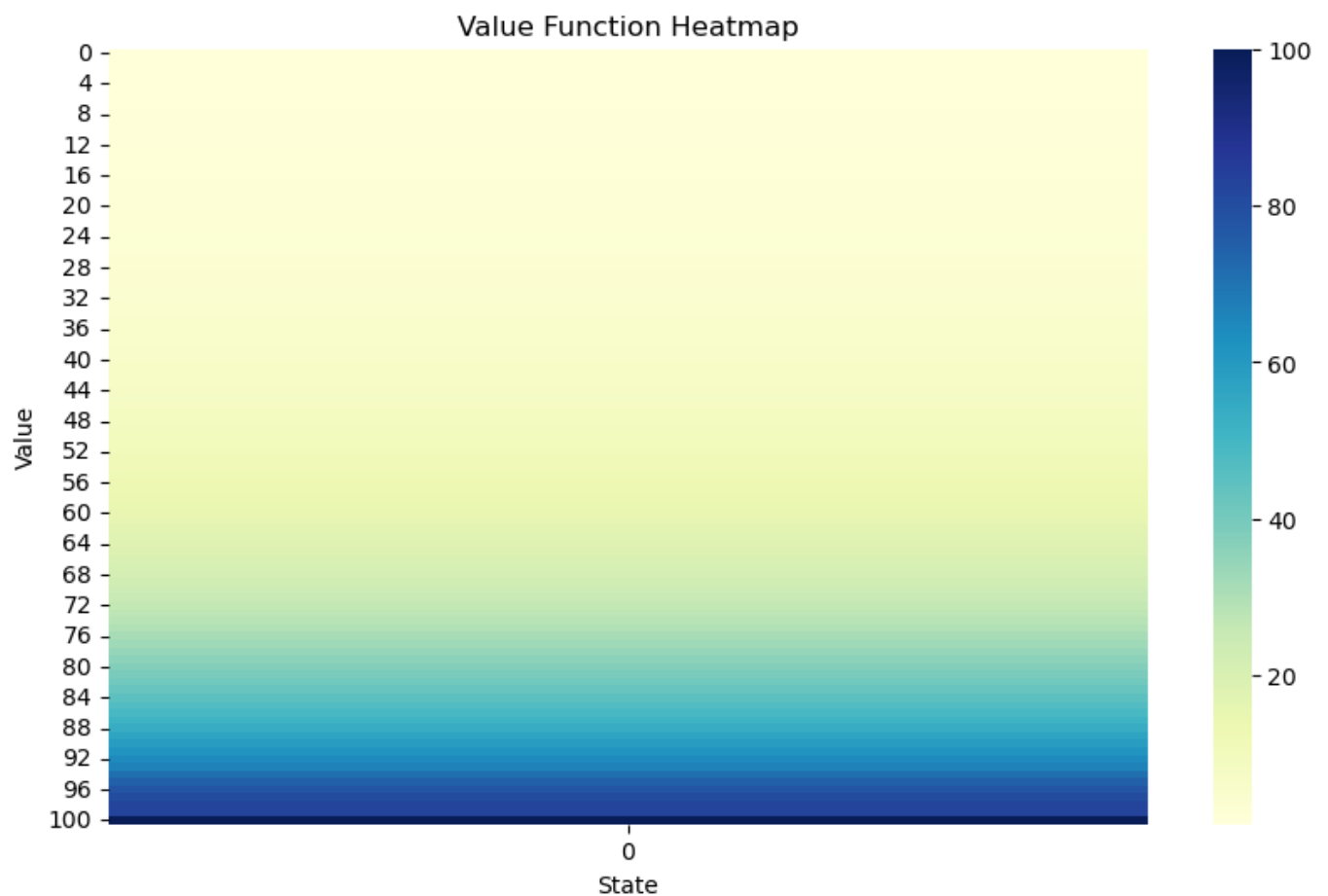
plt.legend()
plt.grid()
plt.show()

Policy Bar Plot
plt.figure(figsize=(10, 6))
plt.bar(range(goal + 1), optimal_policy_vi, label="Optimal Policy", color='orange')
plt.title("Optimal Policy for Value Iteration")
plt.xlabel("State")
plt.ylabel("Action (0 = Stop, 1 = Roll)")
plt.legend()
plt.grid()
plt.show()

Heatmap of Value Function
plt.figure(figsize=(10, 6))
sns.heatmap(optimal_value_vi.reshape(-1, 1), cmap="YlGnBu", cbar=True)
plt.title("Value Function Heatmap")
plt.xlabel("State")
plt.ylabel("Value")
plt.show()

```





## Conclusion (0.5M)

Conclude your assignment in 250 words by discussing the best approach for dice problem with the initial parameters and after changing the parameters.

----write below this line-----

In this assignment, we implemented reinforcement learning agents using dynamic programming

methods, specifically policy iteration and value iteration, to solve the Dice game optimally. The results demonstrate that both methods can effectively determine the optimal policy and value function for reaching a score of exactly 100. Policy iteration iteratively evaluates the value function under the current policy and then improves the policy based on the updated value function until it becomes stable. Value iteration, on the other hand, focuses on iteratively updating the value function directly using the Bellman equation until convergence and then deriving the optimal policy from the value function. When the goal score was changed to 50, the optimal policy and value function adjusted accordingly, illustrating the flexibility of these methods in handling different objectives. Similarly, altering the die to have 8 sides impacted the policy and value function, showing the adaptability of the approach to different game dynamics. Experimenting with different discount factors revealed that the discount rate influences the weight given to future rewards, thereby affecting the optimal policy and value function. A higher discount factor emphasizes future rewards more, leading to different policy decisions compared to a lower discount factor. Visualizing the value function with a heatmap provided insights into how the value of states changes as the player approaches the goal, highlighting the effectiveness of the learned policies. Overall, dynamic programming methods proved to be robust and adaptable for solving the Dice game optimally under various scenarios and parameter settings.