## ⌄  Group No : 43

## Group Member Names:

1. LAKSHMISRINIVAS PERAKAM
2. SHAILESH KUMAR SINGH
3. SUBHRANSU MISHRA
4. JAWAHARLAL RAJAN S

1.**Problem statement**:

- Develop a reinforcement learning agent using dynamic programming to solve the Treasure Hunt problem in a FrozenLake environment. The agent must learn the optimal policy for navigating the lake while avoiding holes and maximizing its treasure collection.

2.**Scenario**:

- A treasure hunter is navigating a slippery 5x5 FrozenLake grid. The objective is to navigate through the lake collecting treasures while avoiding holes and ultimately reaching the exit (goal). Grid positions on a 5x5 map with tiles labeled as S, F, H, G, T. The state includes the current position of the agent and whether treasures have been collected.

## ⌄  Objective

- The agent must learn the optimal policy $\pi^*$ using dynamic programming to maximize its cumulative reward while navigating the lake.

### About the environment

The environment consists of several types of tiles:

- Start (S): The initial position of the agent, safe to step.
- Frozen Tiles (F): Frozen surface, safe to step.
- Hole (H): Falling into a hole ends the game immediately (die, end).
- Goal (G): Exit point; reaching here ends the game successfully (safe, end).
- Treasure Tiles (T): Added to the environment. Stepping on these tiles awards +5 reward but does not end the game.

After stepping on a treasure tile, it becomes a frozen tile (F). The agent earns rewards as follows:

- Reaching the goal (G): +10 reward.
- Falling into a hole (H): -10 reward.
- Collecting a treasure (T): +5 reward.
- Stepping on a frozen tile (F): 0 reward.

### States

- Current position of the agent (row, column).
- A boolean flag (or equivalent) for whether each treasure has been collected.

### Actions

- Four possible moves: up, down, left, right

### Rewards

- Goal (G): +10.
- Treasure (T): +5 per treasure.
- Hole (H): -10.
- Frozen tiles (F): 0.

### Environment

Modify the FrozenLake environment in OpenAI Gym to include treasures (T) at certain positions. Inherit the original FrozenLakeEnv and modify the reset and step methods accordingly. Example grid:


image.png

Double-click (or enter) to edit

**Expected Outcomes:**

1. Create the custom environment by modifying the existing "FrozenLakeNotSlippery-v0" in OpenAI Gym and Implement the dynamic programming using value iteration and policy improvement to learn the optimal policy for the Treasure Hunt problem.
2. Calculate the state-value function (V*) for each state on the map after learning the optimal policy.

3. Compare the agent's performance with and without treasures, discussing the trade-offs in reward maximization.

4. Visualize the agent's direction on the map using the learned policy.

5. Calculate expected total reward over multiple episodes to evaluate performance.

## Import required libraries and Define the custom environment - 2 Marks

```
1 # Import necessary libraries
2 import numpy as np
3 import gym
4 from gym.envs.toy_text.frozen_lake import FrozenLakeEnv
5 import matplotlib.pyplot as plt
6
7 import warnings
8 warnings.filterwarnings("ignore", category=DeprecationWarning)
```

## Custom environment to create the given grid and respective functions that are required for the problem

## Include functions to take an action, get reward, to check if episode is over

```
1 class FrozenLakeTreasureEnv(FrozenLakeEnv):
2     """
3     Custom FrozenLake environment with treasures (T).
4     Inherits from OpenAI Gym's FrozenLakeEnv.
5     """
6     def __init__(self, desc=None, is_slippery=False):
7         """
8         Initializes the environment with a custom grid.
9
10        Args:
11        - desc: Custom description of the grid (list of strings).
12        - is_slippery: If True, makes the environment slippery.
13        """
14        if desc is None:
15            raise ValueError("A custom grid (desc) must be provided for the environment.")
16        super().__init__(desc=desc, map_name=None, is_slippery=is_slippery)
17        self.treasure_positions = [(0, 4), (2, 3), (3, 0)]  # Example treasure locations
18        self.collected_treasures = set()  # To track collected treasures
19    def reset(self):
20        """
21        Resets the environment to its initial state.
22        Clears the list of collected treasures.
23
24        Returns:
25        - The initial state.
26        """
27        self.collected_treasures = set()  # Reset collected treasures
28        return super().reset()
29
30    def step(self, action):
31        """
32        Executes an action in the environment and returns the result.
33        Adds +5 reward for collecting treasures.
34        """
35        # Call the parent class's step method
36        result = super().step(action)
37
38        # Unpack the returned values appropriately
39        if len(result) == 5:
40            next_state, reward, done, info, extra = result
41            # If there's an extra value, you can ignore it or process it based on your needs
42        elif len(result) == 4:
43            next_state, reward, done, info = result
44        elif len(result) == 3:
45            next_state, reward, done = result
46            info = {}  # Default to an empty dictionary if `info` is not returned
47        else:
48            raise ValueError(f"Unexpected number of return values from step: {len(result)}")
49
50        # Calculate the current position of the agent
51        row, col = self.s // self.ncol, self.s % self.ncol  # Convert state index to grid coordinates
52
53        # Check if the agent has stepped on a treasure
```

```
54          if (row, col) in self.treasure_positions and (row, col) not in self.collected_treasures:
55              reward += 5  # Add reward for collecting a treasure
56              self.collected_treasures.add((row, col))  # Mark treasure as collected
57
58          # Check if the episode ends (goal or hole)
59          if self.desc[row, col] in [b'G', b'H']:
60              done = True
61
62          return next_state, reward, done, info
63
64
65      def is_episode_over(self):
66          """
67          Checks if the episode has ended.
68
69          Returns:
70          - True if the episode has ended, False otherwise.
71          """
72          row, col = self.s // self.ncol, self.s % self.ncol  # Current position
73          return self.desc[row][col] in [b'G', b'H']  # End if at Goal or Hole
74
75      def render_custom(self):
76          """
77          Renders the environment with additional info about collected treasures.
78          """
79          print("Environment Grid:")
80          self.render()
81          print(f"Collected Treasures: {len(self.collected_treasures)} / {len(self.treasure_positions)}")
82
```

## ⌄ Value Iteration Algorithm - 1 Mark

```
1 def value_iteration(env, gamma=0.9, theta=1e-4):
2     """
3     Performs value iteration to compute the optimal value function (V*) and policy (π*).
4
5     Parameters:
6     - env: The environment (FrozenLakeTreasureEnv)
7     - gamma: Discount factor
8     - theta: Convergence threshold
9
10    Returns:
11    - V: Optimal value function for all states
12    - policy: Optimal policy (actions for each state)
13    """
14    V = np.zeros(env.observation_space.n)  # Initialize value function
15    policy = np.zeros(env.observation_space.n, dtype=int)  # Initialize policy
16
17    while True:
18        delta = 0  # Tracks the maximum change in the value function
19        for s in range(env.observation_space.n):
20            # Determine if the state is terminal
21            row, col = s // env.ncol, s % env.ncol
22            if env.desc[row, col] in [b'H', b'G']:  # Hole or Goal
23                continue
24
25            # Compute Q-values for all actions
26            q_values = [
27                sum(p * (r + gamma * V[s_]) for p, s_, r, done in env.P[s][a])
28                for a in range(env.action_space.n)
29            ]
30            # Update the value function for state s
31            new_value = max(q_values)
32            delta = max(delta, abs(new_value - V[s]))
33            V[s] = new_value
34            # Update the policy to the action with the highest Q-value
35            policy[s] = np.argmax(q_values)
36
37        # Break if the value function has converged
38        if delta < theta:
39            break
40
41    return V, policy
42
```

## ⌄ Policy Improvement Function - 1 Mark

```
1 def policy_improvement(env, V, gamma=0.9):
2     """
3     Derives an improved policy based on the given value function (V).
4
5     Parameters:
6     - env: The environment (FrozenLakeTreasureEnv)
7     - V: Current value function
8     - gamma: Discount factor
9
10    Returns:
11    - policy: Improved policy
12    """
13    policy = np.zeros(env.observation_space.n, dtype=int)
14    for s in range(env.observation_space.n):
15        if s in env.terminal_states:  # Skip terminal states
16            continue
17        # Compute Q values for all actions
18        q_values = [
19            sum(p * (r + gamma * V[s_]) for p, s_, r, done in env.P[s][a])
20            for a in range(env.action_space.n)
21        ]
22        policy[s] = np.argmax(q_values)  # Choose the best action
23    return policy
24
```

## ∨ Print the Optimal Value Function

```
1 def print_value_function(V, env):
2     """
3     Displays the optimal value function in grid form.
4     """
5     grid = np.array(V).reshape(env.nrow, env.ncol)
6     print("Optimal Value Function:")
7     print(grid)
8
9 custom_desc = [
10    "SFFHT",
11    "FHFFF",
12    "FFFTF",
13    "TFHFF",
14    "FFFFG"
15 ]
16
17 # Assuming env is your environment and V is the optimal value function computed earlier
18 env = FrozenLakeTreasureEnv(desc=custom_desc, is_slippery=False)
19
20 # Compute the optimal value function and policy using value iteration
21 V, policy = value_iteration(env)
22
23 # Print the optimal value function
24 print_value_function(V, env)
25
```

```
Optimal Value Function:
[[0.4782969 0.531441  0.59049   0.        0.729    ]
 [0.531441  0.        0.6561    0.729     0.81     ]
 [0.59049   0.6561    0.729     0.81      0.9      ]
 [0.6561    0.729     0.        0.9       1.       ]
 [0.729     0.81      0.9       1.        0.       ]]
```

## ∨ Visualization of the learned optimal policy - 1 Mark

```
1 def visualize_policy(env, policy):
2     """
3     Visualizes the optimal policy as arrows on the grid.
4     """
5     action_symbols = ['↑', '↓', '←', '→']  # Corresponding to actions 0, 1, 2, 3
6     grid = np.array(env.desc, dtype=str)
7     policy_grid = grid.copy()
8
9     for s in range(env.observation_space.n):
10        row, col = s // env.ncol, s % env.ncol
11        if grid[row, col] in ['H', 'G']:
12            continue
13        policy_grid[row, col] = action_symbols[policy[s]]
14
15    print("Learned Optimal Policy:")
16    for row in policy_grid:
```

```
17         print(' '.join(row))
18
19 custom_desc = [
20     "SFFHT",
21     "FHFFF",
22     "FFFTF",
23     "TFHFF",
24     "FFFFG"
25 ]
26
27 # Updated environment initialization
28 env = FrozenLakeTreasureEnv(desc=custom_desc, is_slippery=False)
29
30 # Compute the optimal value function and policy using value iteration
31 V, policy = value_iteration(env)
32
33 # Visualize the learned optimal policy
34 visualize_policy(env, policy)
```

```
Learned Optimal Policy:
↓ ← ↓ H ↓
↓ H ↓ ↓ ↓
↓ ↓ ← ↓ ↓
↓ ↓ H ↓ ↓
← ← ← ← G
```

## Evaluate the policy - 1 Mark

```
1 def evaluate_policy(env, policy, num_episodes=100):
2     """
3     Evaluates the given policy by running it over multiple episodes.
4
5     Parameters:
6     – env: The environment
7     – policy: The policy to evaluate
8     – num_episodes: Number of episodes to run
9
10    Returns:
11    – mean_reward: Average total reward over all episodes
12    – rewards: List of rewards for each episode
13    """
14    rewards = []
15    for _ in range(num_episodes):
16        state = env.reset()
17        done = False
18        total_reward = 0
19
20        while not done:
21            action = policy[state]
22            state, reward, done, _ = env.step(action)
23            total_reward += reward
24
25        rewards.append(total_reward)
26    mean_reward = np.mean(rewards)
27    return mean_reward, rewards
```

## Main Execution

```
1 if __name__ == "__main__":
2     # Define a custom 5x5 grid
3     custom_desc = [
4         "SFFHT",
5         "FHFFF",
6         "FFFTF",
7         "TFHFF",
8         "FFFFG"
9     ]
10
11    # Initialize the custom environment with the custom grid
12    env = FrozenLakeTreasureEnv(desc=custom_desc, is_slippery=False)
13
14    # Perform value iteration
15    V, policy = value_iteration(env)
16
17    # Print the optimal value function
18    print_value_function(V, env)
19
20    # Visualize the learned policy
21    visualize_policy(env, policy)
```

```
22
23     # Evaluate the policy
24     mean_reward, rewards = evaluate_policy(env, policy)
25     print(f"Mean Reward over {len(rewards)} episodes: {mean_reward}")
26
```

```
Optimal Value Function:
[[0.4782969 0.531441  0.59049   0.        0.729     ]
 [0.531441  0.        0.6561    0.729     0.81      ]
 [0.59049   0.6561    0.729     0.81      0.9       ]
 [0.6561    0.729     0.        0.9       1.        ]
 [0.729     0.81      0.9       1.        0.        ]]
Learned Optimal Policy:
↓ ← ↓ H ↓
↓ H ↓ ↓ ↓
↓ ↓ ← ↓ ↓
↓ ↓ H ↓ ↓
← ← ← ← G
Mean Reward over 100 episodes: 6.0
```