

✓ Alireza Amiri - Project 3

Task 1

```
1 # Import necessary libraries
2 import numpy as np
3 import gym
4 from gym.envs.registration import register
5 import matplotlib.pyplot as plt
6 import pandas as pd
```

✓ 1. Environment Registration

✓ 1.1 Custom Frozen Lake Environment

Register two versions of the Frozen Lake environment: one that is slippery and one that is not. This allows us to compare the performance of the RL algorithms under different conditions.

```
1 # Register custom Frozen Lake environment (non-slippery)
2 register(
3     id='FrozenLakeNotSlippery-v0',
4     entry_point='gym.envs.toy_text:FrozenLakeEnv',
5     kwargs={'desc': [
6         "SFFFFFFH",
7         "FFHHHHFH",
8         "FFHHFFFF",
9         "FFFHFHGG",
10        "FFHHFHFF",
11        "FFHFFHFF",
12        "FHFFFFFF",
13        "FFFFFFF",
14    ], 'is_slippery': False}
15 )
16
17 # Register custom Frozen Lake environment (slippery)
18 register(
19     id='FrozenLakeSlippery-v0',
20     entry_point='gym.envs.toy_text:FrozenLakeEnv',
21     kwargs={'desc': [
22         "SFFFFFFH",
23         "FFHHHHFH",
24         "FFHHFFFF",
25         "FFFHFHGG",
26         "FFHHFHFF",
27         "FFHFFHFF",
28         "FHFFFFFF",
29         "FFFFFFF",
30     ], 'is_slippery': True}
31 )
32
33 # Create environments
34 env_not_slippery = gym.make('FrozenLakeNotSlippery-v0', new_step_api=True)
35 env_slippery = gym.make('FrozenLakeSlippery-v0', new_step_api=True)
36
```

✓ 2. Q-Learning and SARSA Algorithms

✓ 2.1 Algorithm Parameters

Define the parameters for the Q-Learning and SARSA algorithms.

```
1 # Algorithm parameters
2 alpha = 0.8
3 epsilon = 1.0
4 epsilon_decay = 0.9999999999
5 min_epsilon = 0.01
6 episodes = 100000
7 max_steps = 100
8 gamma_values = [0.5, 0.7, 0.99]
9
```

✓ 2.2 Q-Learning Implementation

Q-Learning is a value-based method of reinforcement learning. It updates the Q-value of the current state-action pair using the Bellman equation.

```
1 def q_learning(env, gamma):
2     Q = np.zeros((env.observation_space.n, env.action_space.n))
3     rewards = []
4     epsilons = []
5     epsilon = 1.0
6
7     def choose_action(state, epsilon, visited_states):
8         if np.random.rand() < epsilon or state in visited_states:
9             return env.action_space.sample()
10        else:
11            return np.argmax(Q[state, :])
12
13    for episode in range(episodes):
14        if episode % 1000 == 0:
15            print(f"Q-Learning: Episode {episode}")
16        state = env.reset()
17        state = state[0] if isinstance(state, tuple) else state
18        total_reward = 0
19        visited_states = set()
20
21        for step in range(max_steps):
22            visited_states.add(state)
23            action = choose_action(state, epsilon, visited_states)
24            next_state, reward, done, truncated, info = env.step(action)
25
26            Q[state, action] = Q[state, action] + alpha * (reward + gamma * np.max(Q[next_state, :]) - Q[state, action])
27
28            state = next_state
29            total_reward += reward
30
31            if done:
32                break
33
34        epsilon = max(min_epsilon, epsilon * epsilon_decay)
35        rewards.append(total_reward)
36        epsilons.append(epsilon)
37
38    return Q, rewards, epsilons
39
```

✓ 2.3 SARSA Implementation

SARSA (State-Action-Reward-State-Action) is an on-policy method. It updates the Q-value using the action taken in the next state.

```
1 def sarsa(env, gamma):
2     Q = np.zeros((env.observation_space.n, env.action_space.n))
3     rewards = []
4     epsilons = []
5     epsilon = 1.0
6
7     def choose_action(state, epsilon):
8         if np.random.rand() < epsilon:
9             return env.action_space.sample()
10        else:
11            return np.argmax(Q[state, :])
12
13    for episode in range(episodes):
14        if episode % 1000 == 0:
15            print(f"SARSA: Episode {episode}")
16        state = env.reset()
17        state = state[0] if isinstance(state, tuple) else state
18        total_reward = 0
19
20        action = choose_action(state, epsilon)
21
22        for step in range(max_steps):
23            next_state, reward, done, truncated, info = env.step(action)
24            next_state = next_state if isinstance(next_state, int) else next_state[0]
25            next_action = choose_action(next_state, epsilon)
26
27            Q[state, action] = Q[state, action] + alpha * (reward + gamma * Q[next_state, next_action] - Q[state, action])
28
29            state = next_state
30            action = next_action
31            total_reward += reward
32
33            if done:
34                break
35
36        epsilon = max(min_epsilon, epsilon * epsilon_decay)
37        rewards.append(total_reward)
38        epsilons.append(epsilon)
39
40    return Q, rewards, epsilons
41
```

✓ 3. Display and Analysis

✓ 3.1 Display Optimal Policy

A function to display the optimal policy derived from the Q-values.

```
1 def display_policy(Q):
2     actions = ["←", "↓", "→", "↑"]
3     optimal_policy = np.argmax(Q, axis=1).reshape(8, 8)
4     policy_display = np.vectorize(lambda x: actions[x])(optimal_policy)
5     for row in policy_display:
6         print(" ".join(row))
7
```

✓ 3.2 Train and Compare Algorithms

Train both Q-Learning and SARSA algorithms with different gamma values and compare their performance.

```

1 # Train Q-Learning and SARSA with different gamma values
2 results_q_learning_not_slippery = []
3 results_q_learning_slippery = []
4 results_sarsa_not_slippery = []
5 results_sarsa_slippery = []
6
7 for gamma in gamma_values:
8     print(f"Training Q-Learning (slippery=False) with gamma={gamma}")
9     Q, rewards, epsilons = q_learning(env_not_slippery, gamma)
10    results_q_learning_not_slippery.append({'gamma': gamma, 'Q': Q, 'rewards': rewards, 'epsilons': epsilons})
11
12    print(f"Training Q-Learning (slippery=True) with gamma={gamma}")
13    Q, rewards, epsilons = q_learning(env_slippery, gamma)
14    results_q_learning_slippery.append({'gamma': gamma, 'Q': Q, 'rewards': rewards, 'epsilons': epsilons})
15
16    print(f"Training SARSA (slippery=False) with gamma={gamma}")
17    Q, rewards, epsilons = sarsa(env_not_slippery, gamma)
18    results_sarsa_not_slippery.append({'gamma': gamma, 'Q': Q, 'rewards': rewards, 'epsilons': epsilons})
19
20    print(f"Training SARSA (slippery=True) with gamma={gamma}")
21    Q, rewards, epsilons = sarsa(env_slippery, gamma)
22    results_sarsa_slippery.append({'gamma': gamma, 'Q': Q, 'rewards': rewards, 'epsilons': epsilons})
23

```

 [Show hidden output](#)

3.3 Plot Results

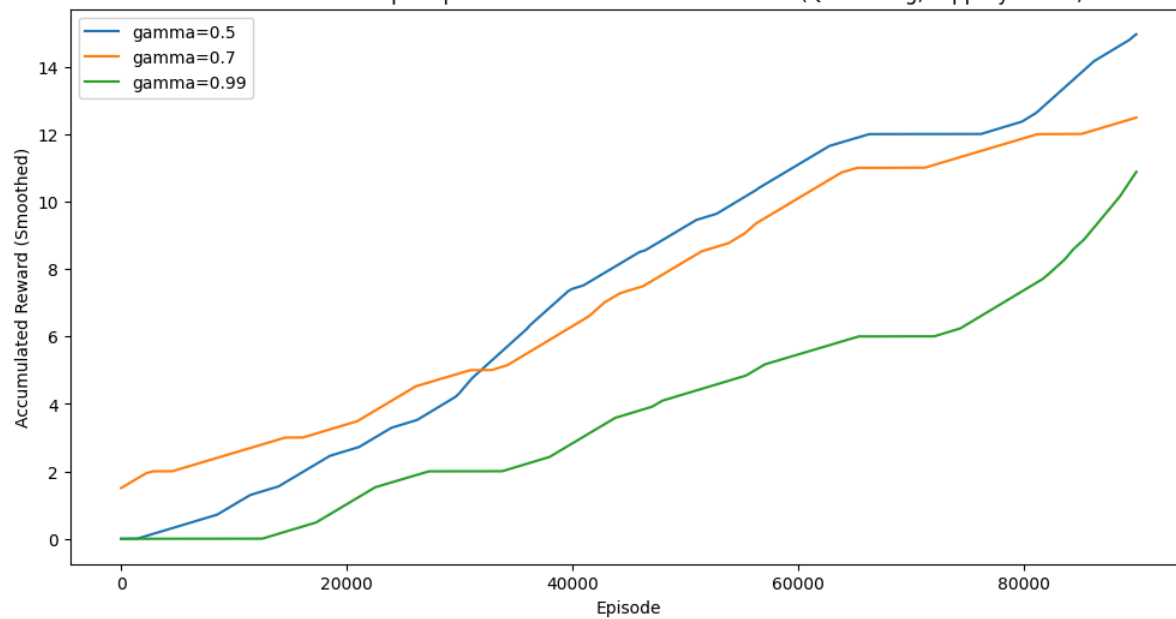
Plot the accumulated rewards for each algorithm.

```

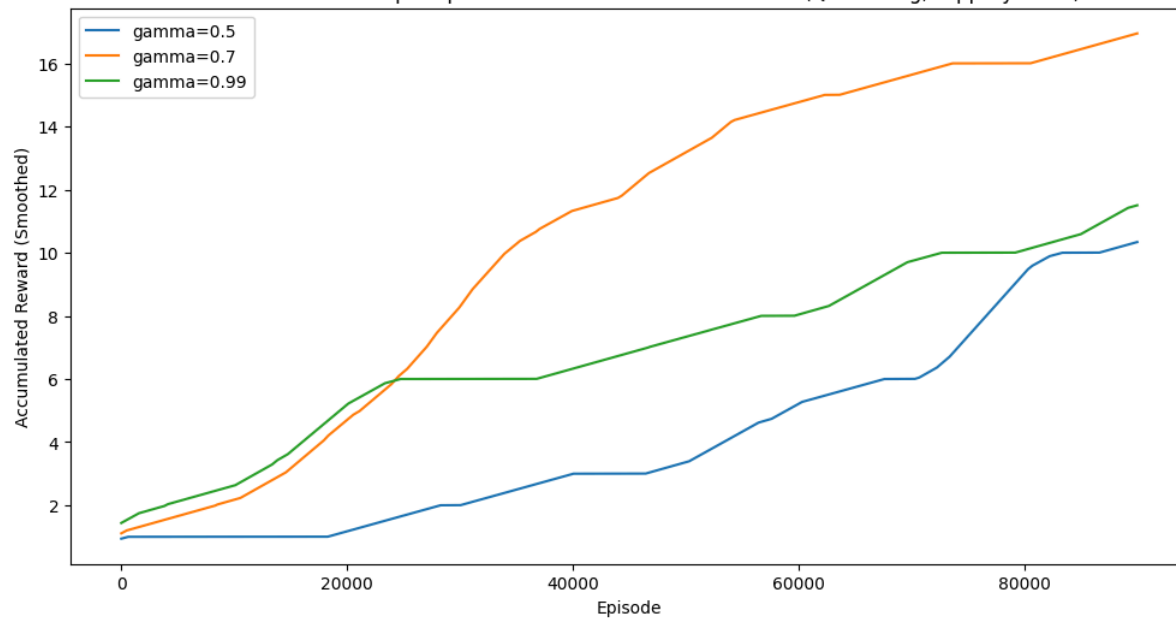
1 def plot_rewards_smooth(results, title, filename, window_size=100):
2     plt.figure(figsize=(12, 6))
3     for result in results:
4         accumulated_rewards = np.cumsum(result['rewards'])
5         smoothed_rewards = moving_average(accumulated_rewards, window_size)
6         plt.plot(smoothed_rewards, label=f'gamma={result["gamma"]}')
7     plt.xlabel('Episode')
8     plt.ylabel('Accumulated Reward (Smoothed)')
9     plt.title(title)
10    plt.legend()
11    plt.savefig(filename)
12    plt.show()
13
14 # Call the function with the desired window size
15 plot_rewards_smooth(results_q_learning_not_slippery, 'Accumulated Reward per Episode for Different Gamma Values (Q-Learning, slippery=False)', 'q_learning_not_slippery.png', window_size=10000)
16 plot_rewards_smooth(results_q_learning_slippery, 'Accumulated Reward per Episode for Different Gamma Values (Q-Learning, slippery=True)', 'q_learning_slippery.png', window_size=10000)
17 plot_rewards_smooth(results_sarsa_not_slippery, 'Accumulated Reward per Episode for Different Gamma Values (SARSA, slippery=False)', 'sarsa_not_slippery.png', window_size=10000)
18 plot_rewards_smooth(results_sarsa_slippery, 'Accumulated Reward per Episode for Different Gamma Values (SARSA, slippery=True)', 'sarsa_slippery.png', window_size=10000)
19

```

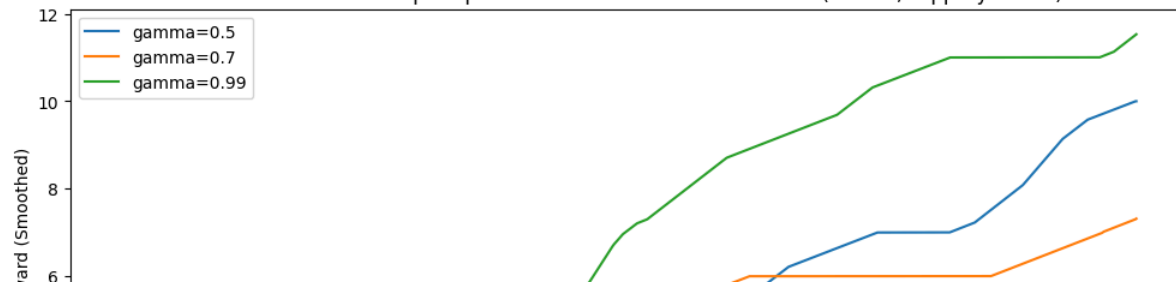
Accumulated Reward per Episode for Different Gamma Values (Q-Learning, slippery=False)

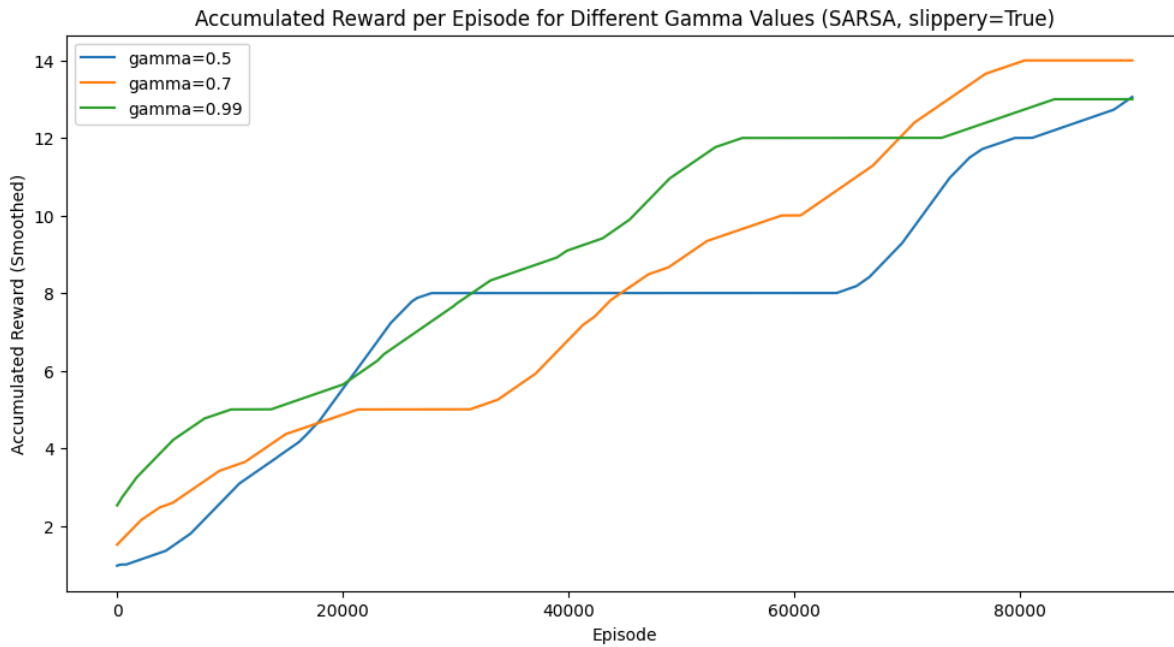
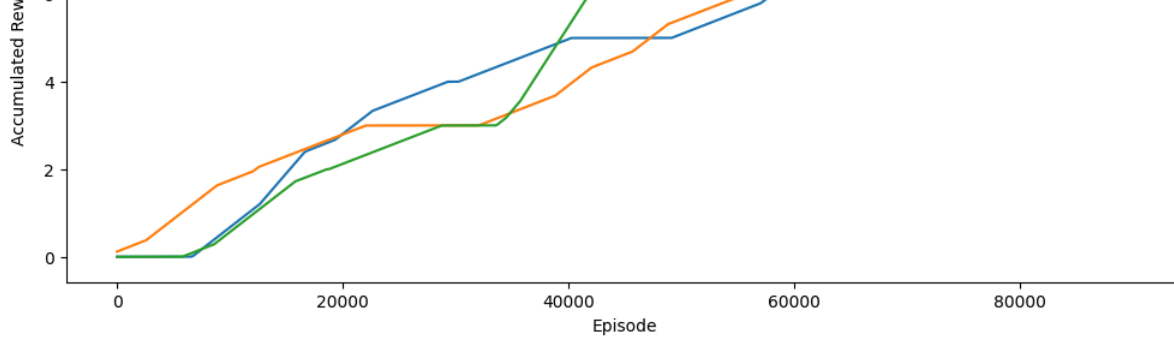


Accumulated Reward per Episode for Different Gamma Values (Q-Learning, slippery=True)



Accumulated Reward per Episode for Different Gamma Values (SARSA, slippery=False)





3.4 Display Optimal Policies

Display the optimal policies for both Q-Learning and SARSA in slippery and non-slippery environments.

```
1 # Display optimal policies
2 print("Optimal Policy for Q-Learning (slippery=False):")
3 for result in results_q_learning_not_slippery:
4     print(f"gamma={result['gamma']}")
5     display_policy(result['Q'])
6
7 print("Optimal Policy for Q-Learning (slippery=True):")
8 for result in results_q_learning_slippery:
9     print(f"gamma={result['gamma']}")
10    display_policy(result['Q'])
11
12 print("Optimal Policy for SARSA (slippery=False):")
13 for result in results_sarsa_not_slippery:
14     print(f"gamma={result['gamma']}")
15     display_policy(result['Q'])
16
17 print("Optimal Policy for SARSA (slippery=True):")
18 for result in results_sarsa_slippery:
19     print(f"gamma={result['gamma']}")
```

```
20     display_policy(result['Q'])
21
↩
↑ ↓ ↑ ↑ ↓ → → ←
→ ← ← ← ← → ←
→ ↑ ← ← ← ↓ ↑ ↓
↑ ↑ ← ← ↓ ← ← ←
← ← ← → ← ← ↑
← ← ← ↑ ← ← ↓ ←
→ ← ← ↑ → ↑ ↑ ↑
↓ → → → → → →
gamma=0.7
↑ → ↑ ↓ ↑ ← → ←
← ↑ ← ← ← ← ↑ ←
← → ← ← ↓ → → ←
← ↑ ← ← ← ← ← ←
→ → ← ← → ← ← →
↑ ↑ ← ↓ ← ← ↑ ↑
← ← → → ↑ → → ←
↑ ↑ ↑ → → ← ← ←
gamma=0.99
↑ → ↓ → ↑ ↑ ← ←
→ ← ← ← ← → ←
↑ ↑ ← ← ↓ → → ↓
← ← ← ← ← ← ← ←
← ← ← ← ↓ ← ← →
← ← ← → ← ← → →
← ← → ↑ → ↓ ↑ →
↓ ↓ ↑ → → ↑ →
Optimal Policy for SARSA (slippery=False):
gamma=0.5
↑ ← → → → ↑ ↓ ←
↑ ↑ ← ← ← ↓ ← ←
↑ ↑ ← ← → → ↓
↑ ← ← ← ← ← ← ←
↑ ← ← ← ↓ ← ← ↑
← ↑ ← ↓ ← ← ↓ ↑
↑ ← ↓ ↓ ← ↓ ← ↑
↓ ← → ↓ ← ↓ ← ↓
gamma=0.7
↑ ← → → ↑ → ↑ ←
← ↑ ← ← ← ↓ ← ←
↑ ↑ ← ← → → ↓
↓ ↑ ← ← ← ← ← ←
↑ ↑ ← ← ↓ ← ← ↑
← ← ← ↓ ← → ↓
↑ ← → ↓ → → ↑ ↑
→ → → ↓ ↑ ← ← ←
gamma=0.99
↑ ↑ ← → → → ↑ ←
↑ ↓ ← ← ← ↓ ← ←
→ ← ← ← ← → ↓
→ ← ← ← ← ← ← ←
← ↑ ← ← ← ← ↑
↑ ↑ ← ↓ ↓ ← ← ↓
↓ ← → → → ↓ ← ←
↑ ← ← ↓ → ↓ ← ↑
Optimal Policy for SARSA (slippery=True):
gamma=0.5
↓ → ↑ ↑ ↑ ↑ ↑ ←
↓ ↑ ← ← ← ← ← ←
↑ → ← ← ↑ ↑ → →
```

Double-click (or enter) to edit

1 Start coding or [generate](#) with AI.