# CSCI B659: Reinforcement Learning
# Assignment 2: Everything Printout

LJ Huang

Spring 2025

## Contents

# 1 Task 1: SARSA

In this section, we present our experimental results for Task 1, where we applied the SARSA algorithm to three distinct environments: FrozenLake4, FrozenLake8, and CartPole. Our goal is to demonstrate how SARSA performs across different reinforcement learning problems – from discrete navigation tasks in FrozenLake to the more dynamic control challenge of balancing a pole in CartPole.

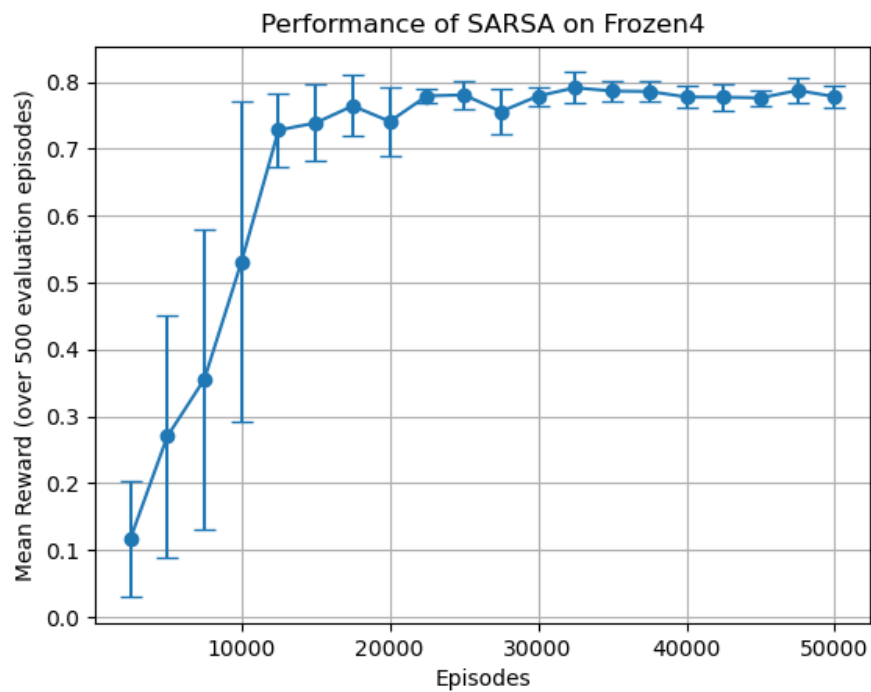## 1.1 Experimental Plots for Task 1

### 1.1.1 SARSA FrozenLake4



Figure 1: SARSA performance on FrozenLake4

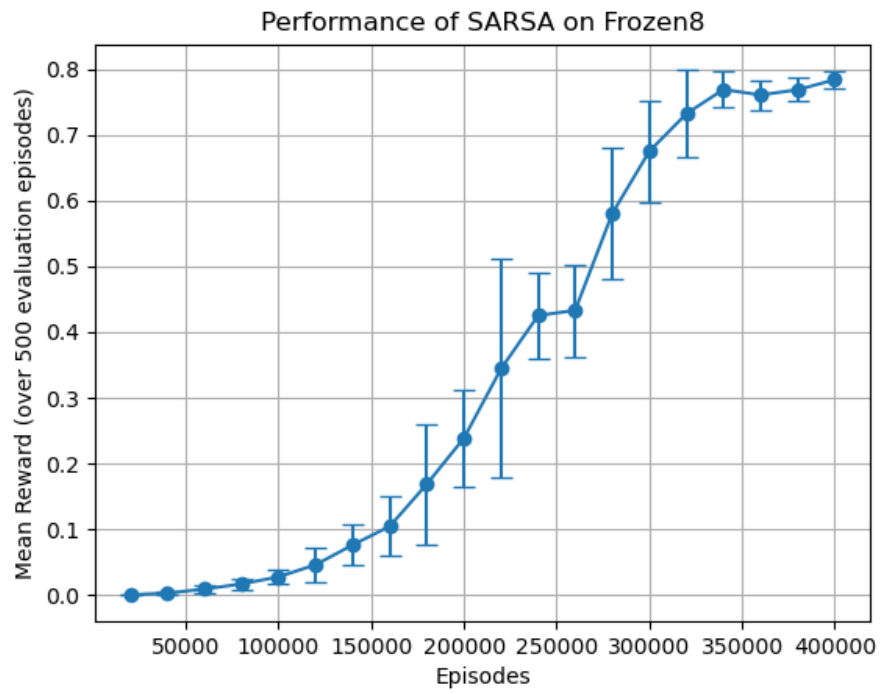### 1.1.2   SARSA FrozenLake8



Figure 2: SARSA performance on FrozenLake8
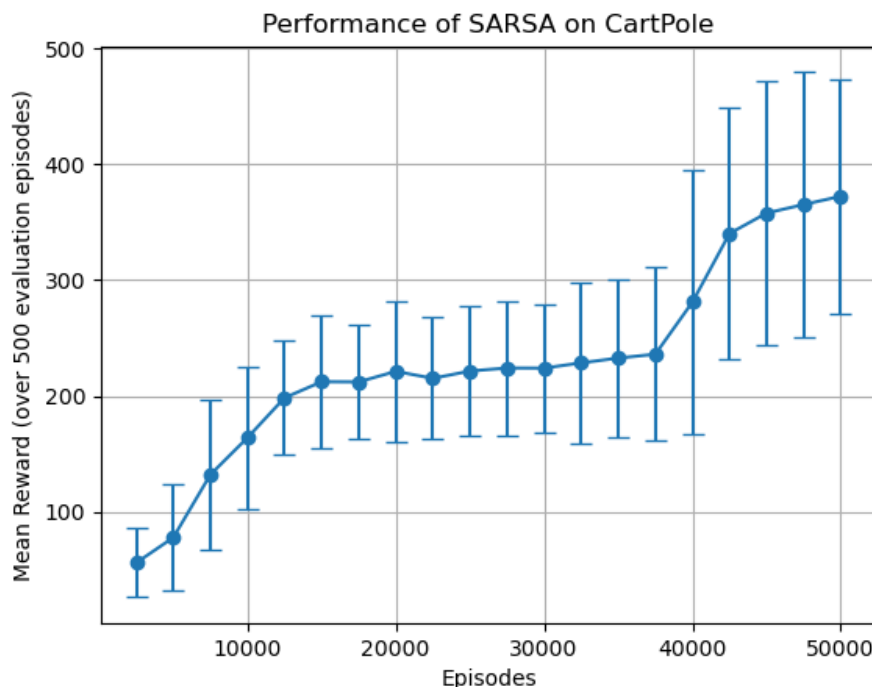
### 1.1.3 SARSA CartPole



Figure 3: SARSA performance on CartPole

## 1.2 Observations on Task 1

**FrozenLake4:** SARSA converges within about 10,000 episodes to an average reward near 0.8, with error bars indicating stable performance after roughly 20,000–30,000 episodes.

**FrozenLake8:** Due to its larger state space, the agent requires hundreds of thousands of episodes to stabilize, eventually achieving mean rewards between 0.8 and 0.85 after extensive exploration.

**CartPole:** The mean reward steadily increases during the first 30,000 episodes and then accelerates, reaching values near 400–450 by 50,000 episodes, which reflects improved balance over time despite some variability.

**Variant of SARSA Used:** We employed an on-policy SARSA algorithm with an $\epsilon$-greedy exploration strategy, where $\epsilon$ decays linearly from 0.5 to 0. Key hyperparameters are:

- Learning rate: $\alpha = 0.01$

- Discount factor: $\gamma = 0.999$

**Sensitivity to Hyperparameters:** SARSA is sensitive to these settings. A too-large $\alpha$ may cause unstable updates, while a too-small $\alpha$ can slow convergence. Similarly, if $\epsilon$ decays too quickly, the agent may exploit prematurely; if it remains high too long, convergence is delayed.

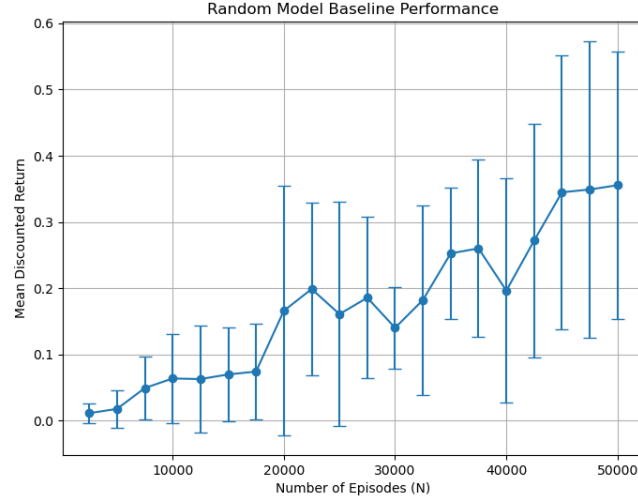# 2  Task 2: Model-Based RL

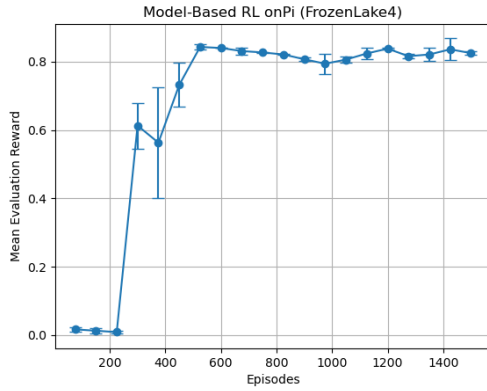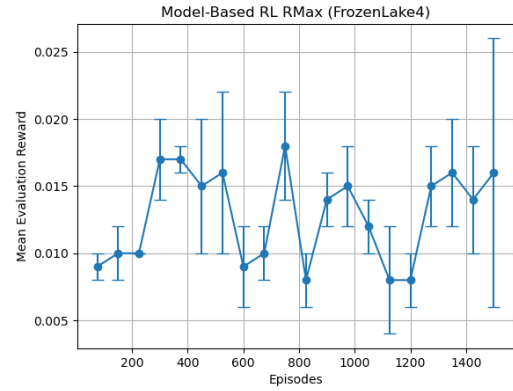## 2.1  Experimental Plots for Task 2



Figure 4: Random Model Baseline (FrozenLake)



(a) onPi Variant on FrozenLake4

(b) RMax Variant on FrozenLake4

Figure 5: Results on FrozenLake4

(a) onPi Variant on FrozenLake8



(b) RMax Variant on FrozenLake8

Figure 6: Results on FrozenLake8

## 2.2 Observations on Task 2

**Success of Each Algorithm:**

- **Random:** Learns a workable policy eventually but is slow and requires large amounts of data.

- **onPi:** Achieves good performance on FrozenLake4 quickly and learns a reasonable policy on FrozenLake8, indicating faster convergence.

- **RMax:** RMax did not perform well in our experiments, possibly due to code implementation issues in `Model_Based_RL.py` or the need for additional parameter tuning.

**Performance and Speed of Learning:**

- **onPi** demonstrates the fastest and most stable learning, particularly in FrozenLake4.

- **Random** improves only with large data collection and is slower overall.

- **RMax** may have potential when correctly implemented and tuned but did not show a clear advantage in these experiments.

# 3  Code Snapshot

Snapshot of codes provided in this section used in this assignment.

## 3.1  SARSA Implementation

```python
import numpy as np
import matplotlib.pyplot as plt
from pp2starter import prepCartPole, prepFrozen4, prepFrozen8


def choose_action(Q, state, nA, epsilon):
    if np.random.rand() < epsilon:
        return np.random.choice(nA)
    else:
        best_actions = np.flatnonzero(Q[state] == np.max(Q[state]))
        return np.random.choice(best_actions)

def sarsa(env, phi, nS, nA, episodes, max_steps, alpha, gamma, initial_epsilon
    =0.5, initial_Q_value=0.0):
    Q = np.full((nS, nA), initial_Q_value, dtype=np.float64)

    for ep in range(episodes):
        epsilon = initial_epsilon * (episodes - ep - 1) / (episodes - 1)

        observation, info = env.reset()
        state = phi(observation)
        action = choose_action(Q, state, nA, epsilon)

        for t in range(max_steps):
            observation, reward, terminated, truncated, info = env.step(action)
            next_state = phi(observation)
            if terminated or truncated:
                delta = reward - Q[state, action]
                Q[state, action] += alpha * delta
                break
            else:
                next_action = choose_action(Q, next_state, nA, epsilon)
                delta = reward + gamma * Q[next_state, next_action] - Q[state,
                    action]
                Q[state, action] += alpha * delta
                state, action = next_state, next_action
    return Q

def evaluate_policy(env, phi, Q, num_episodes=500):
    rewards = []
    for _ in range(num_episodes):
        observation, info = env.reset()
        state = phi(observation)
        total_reward = 0
        for _ in range(env._max_episode_steps):
            best_actions = np.flatnonzero(Q[state] == np.max(Q[state]))
            action = np.random.choice(best_actions)
            observation, reward, terminated, truncated, info = env.step(action)
            total_reward += reward
            state = phi(observation)
            if terminated or truncated:
                break
```

```
51          rewards.append(total_reward)
52      return np.mean(rewards)
53
54 def run_experiments(prep_fn, total_episodes, eval_interval, n_repeats, alpha,
       gamma,
55                      initial_epsilon=0.5, initial_Q_value=0.0):
56      nS, nA, env, phi, dname = prep_fn()
57      max_steps = env._max_episode_steps
58      eval_points = list(range(eval_interval, total_episodes + 1, eval_interval))
59      all_scores = np.zeros((n_repeats, len(eval_points)))
60
61      for rep in range(n_repeats):
62          print(f"Starting repeat {rep+1}/{n_repeats}...")
63          Q = np.full((nS, nA), initial_Q_value, dtype=np.float64)
64          eval_counter = 0
65          for ep in range(total_episodes):
66              epsilon = initial_epsilon * (total_episodes - ep - 1) / (
                    total_episodes - 1)
67              observation, info = env.reset()
68              state = phi(observation)
69              action = choose_action(Q, state, nA, epsilon)
70
71              for t in range(max_steps):
72                  observation, reward, terminated, truncated, info = env.step(action
                        )
73                  next_state = phi(observation)
74                  if terminated or truncated:
75                      delta = reward - Q[state, action]
76                      Q[state, action] += alpha * delta
77                      break
78                  else:
79                      next_action = choose_action(Q, next_state, nA, epsilon)
80                      delta = reward + gamma * Q[next_state, next_action] - Q[state,
                            action]
81                      Q[state, action] += alpha * delta
82                      state, action = next_state, next_action
83
84              if (ep + 1) % eval_interval == 0:
85                  score = evaluate_policy(env, phi, Q)
86                  all_scores[rep, eval_counter] = score
87                  eval_counter += 1
88
89          env.close()
90
91      means = np.mean(all_scores, axis=0)
92      stds = np.std(all_scores, axis=0)
93      return dname, eval_points, means, stds
94
95 def plot_results(dname, eval_points, means, stds):
96      plt.figure()
97      plt.errorbar(eval_points, means, yerr=stds, fmt='-o', capsize=5)
98      plt.title(f'Performance of SARSA on {dname}')
99      plt.xlabel('Episodes')
100     plt.ylabel('Mean Reward (over 500 evaluation episodes)')
101     plt.grid(True)
102     plt.show()
103
104
105 if __name__ == '__main__':
```

```python
106     alpha = 0.01
107     gamma = 0.999
108     n_repeats = 10
109     initial_epsilon = 0.5
110
111     total_episodes_f4 = 50000
112     eval_interval_f4 = total_episodes_f4 // 20
113     dname, eval_points, means, stds = run_experiments(prepFrozen4,
            total_episodes_f4, eval_interval_f4,
114                                             n_repeats, alpha, gamma,
115                                             initial_epsilon=
                                                    initial_epsilon,
116                                             initial_Q_value=0.0)
117     plot_results(dname, eval_points, means, stds)
118
119     total_episodes_f8 = 400000
120     eval_interval_f8 = total_episodes_f8 // 20
121     dname, eval_points, means, stds = run_experiments(prepFrozen8,
            total_episodes_f8, eval_interval_f8,
122                                             n_repeats, alpha, gamma,
123                                             initial_epsilon=
                                                    initial_epsilon,
124                                             initial_Q_value=0.0)
125     plot_results(dname, eval_points, means, stds)
126
127     total_episodes_cp = 50000
128     eval_interval_cp = total_episodes_cp // 20
129     dname, eval_points, means, stds = run_experiments(prepCartPole,
            total_episodes_cp, eval_interval_cp,
130                                             n_repeats, alpha, gamma,
131                                             initial_epsilon=
                                                    initial_epsilon,
132                                             initial_Q_value=0.0)
133     plot_results(dname, eval_points, means, stds)
```

Listing 1: SARSA.py

## 3.2   Model-Based RL Implementation

```python
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from pp1starter import prepFrozen
4 from VI_PI_MPI import value_iteration, policy_evaluation
5 from pp2starter import prepFrozen4, prepFrozen8
6
7 def collect_data(env, N):
8     counts = {}
9     transition_counts = {}
10    reward_sums = {}
11
12    for episode in range(N):
13        observation, info = env.reset()
14        terminated = False
15        truncated = False
16        while not (terminated or truncated):
17            action = env.action_space.sample()
18            next_obs, reward, terminated, truncated, info = env.step(action)
```

```python
                key = (observation, action)
                counts[key] = counts.get(key, 0) + 1
                transition_key = (observation, action, next_obs)
                transition_counts[transition_key] = transition_counts.get(
                    transition_key, 0) + 1
                reward_sums[key] = reward_sums.get(key, 0) + reward
                observation = next_obs
    return counts, transition_counts, reward_sums

def estimate_model(counts, transition_counts, reward_sums, nS, nA):
    P_est = []
    for s in range(nS):
        action_list = []
        for a in range(nA):
            action_list.append(None)
        P_est.append(action_list)

    for s in range(nS):
        for a in range(nA):
            key = (s, a)
            total = counts.get(key, 0)
            outcomes = []
            if total > 0:
                s_primes = set()
                for s0, a0, s_ in transition_counts.keys():
                    if s0 == s and a0 == a:
                        s_primes.add(s_)
                for s_prime in s_primes:
                    transition_key = (s, a, s_prime)
                    count_sas = transition_counts.get(transition_key, 0)
                    p_est = count_sas / total
                    r_est = reward_sums.get(key, 0) / total
                    outcomes.append((p_est, s_prime, r_est, False))
                P_est[s][a] = outcomes
            else:
                P_est[s][a] = [(1.0, s, 0.0, False)]
    return P_est

def run_random_experiment(env, nS, nA, discount, tolerance, N_values,
    num_experiments=10):
    results = {N: [] for N in N_values}
    for N in N_values:
        print(f"[Random] Data collection: N = {N}")
        for exp in range(num_experiments):
            counts, transition_counts, reward_sums = collect_data(env, N)
            P_est = estimate_model(counts, transition_counts, reward_sums, nS, nA)
            policy_est, _, _, _ = value_iteration(P_est, nS, nA, discount,
                tolerance, max_iter=500, env=None)
            score = policy_evaluation(env, policy_est, discount, episodes=500)
            results[N].append(score)
            print(f"  Experiment {exp+1}: Score = {score:.4f}")
    return results


def choose_action(Q, state, nA, epsilon):
    if np.random.rand() < epsilon:
        return np.random.choice(nA)
    else:
        best_actions = np.flatnonzero(Q[state] == np.max(Q[state]))
```

```python
 75            return np.random.choice(best_actions)
 76
 77 def build_model(data, nS, nA, variant='standard', M=10, R=1):
 78     counts = np.zeros((nS, nA, nS))
 79     rewards = np.zeros((nS, nA, nS))
 80     sa_counts = np.zeros((nS, nA))
 81     dead_end = np.zeros(nS, dtype=bool)
 82
 83     for (s, a, r, s_next, terminated) in data:
 84         if s < nS and s_next < nS:
 85             counts[s, a, s_next] += 1
 86             rewards[s, a, s_next] += r
 87             sa_counts[s, a] += 1
 88             if terminated:
 89                 dead_end[s_next] = True
 90
 91     if variant == 'standard' or variant == 'onPi':
 92         P = np.zeros((nS, nA, nS))
 93         R_model = np.zeros((nS, nA))
 94         for s in range(nS):
 95             for a in range(nA):
 96                 if sa_counts[s, a] > 0:
 97                     P[s, a, :] = counts[s, a, :] / sa_counts[s, a]
 98                     R_model[s, a] = np.sum(rewards[s, a, :]) / sa_counts[s, a]
 99                 else:
100                     P[s, a, s] = 1.0
101                     R_model[s, a] = 0.0
102         return P, R_model
103     elif variant == 'RMax':
104         nS_new = nS + 1
105         P = np.zeros((nS_new, nA, nS_new))
106         R_model = np.zeros((nS_new, nA))
107         for s in range(nS):
108             for a in range(nA):
109                 if sa_counts[s, a] >= M:
110                     P[s, a, :nS] = counts[s, a, :] / sa_counts[s, a]
111                     R_model[s, a] = np.sum(rewards[s, a, :]) / sa_counts[s, a]
112                 else:
113                     if not dead_end[s]:
114                         P[s, a, nS] = 1.0
115                         R_model[s, a] = R
116                     else:
117                         P[s, a, s] = 1.0
118                         R_model[s, a] = 0.0
119         for a in range(nA):
120             P[nS, a, nS] = 1.0
121             R_model[nS, a] = 0.0
122         return P, R_model
123
124 def value_iteration_model(P, R_model, gamma, threshold=1e-6, max_iter=10000):
125     nS, nA, _ = P.shape
126     V = np.zeros(nS)
127     for _ in range(max_iter):
128         V_prev = V.copy()
129         for s in range(nS):
130             Q_s = np.zeros(nA)
131             for a in range(nA):
132                 Q_s[a] = R_model[s, a] + gamma * np.sum(P[s, a, :] * V_prev)
133             V[s] = np.max(Q_s)
```

```python
134            if np.max(np.abs(V - V_prev)) < threshold:
135                break
136        Q = np.zeros((nS, nA))
137        for s in range(nS):
138            for a in range(nA):
139                Q[s, a] = R_model[s, a] + gamma * np.sum(P[s, a, :] * V)
140        return V, Q

142    def evaluate_policy_model(env, phi, Q, num_episodes=500):
143        rewards_eval = []
144        for _ in range(num_episodes):
145            observation, info = env.reset()
146            state = phi(observation)
147            total_reward = 0
148            done = False
149            while not done:
150                best_actions = np.flatnonzero(Q[state] == np.max(Q[state]))
151                action = np.random.choice(best_actions)
152                observation, reward, terminated, truncated, info = env.step(action)
153                total_reward += reward
154                state = phi(observation)
155                done = terminated or truncated
156            rewards_eval.append(total_reward)
157        return np.mean(rewards_eval)

159    def model_based_rl(prep_fn, total_episodes, eval_interval, n_repeats,
160                       gamma=0.999, initial_epsilon=0.5, variant='onPi', M=10, R=1):
161        nS, nA, env, phi, dname = prep_fn()
162        quality_all = np.zeros((n_repeats, total_episodes // eval_interval))

164        for rep in range(n_repeats):
165            print(f"[{variant}] Starting repeat {rep+1}/{n_repeats}...")
166            data = []
167            Q = np.zeros((nS, nA))
168            eval_counter = 0
169            for ep in range(total_episodes):
170                if (ep + 1) % eval_interval == 0:
171                    P_eval, R_eval = build_model(data, nS, nA, variant='standard')
172                    _, Q_eval = value_iteration_model(P_eval, R_eval, gamma)
173                    quality = evaluate_policy_model(env, phi, Q_eval)
174                    quality_all[rep, eval_counter] = quality
175                    eval_counter += 1
176                    print(f"  Episode {ep+1}/{total_episodes}: Eval quality = {quality
                        :.2f}")
177                observation, info = env.reset()
178                state = phi(observation)
179                done = False
180                while not done:
181                    if variant == 'onPi':
182                        action = choose_action(Q, state, nA, epsilon=initial_epsilon)
183                    elif variant == 'RMax':
184                        action = np.argmax(Q[state])
185                    else:
186                        raise ValueError("Unknown variant")
187                    observation, reward, terminated, truncated, info = env.step(action
                        )
188                    next_state = phi(observation)
189                    data.append((state, action, reward, next_state, terminated))
190                    state = next_state
```

```python
191                    done = terminated or truncated
192            env.close()
193        return dname, quality_all
194
195 def plot_results(eval_points, quality_all, title):
196        mean_quality = np.mean(quality_all, axis=0)
197        std_quality = np.std(quality_all, axis=0)
198        plt.figure()
199        plt.errorbar(eval_points, mean_quality, yerr=std_quality, fmt='-o', capsize=5)
200        plt.title(title)
201        plt.xlabel('Episodes')
202        plt.ylabel('Mean Evaluation Reward')
203        plt.grid(True)
204        plt.show()
205
206
207 if __name__ == '__main__':
208
209        env_random, P_random, nS, nA, dname_random = prepFrozen()
210        tolerance = 0.001
211        discount = 0.999
212        N_values = range(2500, 50001, 2500)
213        random_results = run_random_experiment(env_random, nS, nA, discount, tolerance
                , N_values, num_experiments=10)
214
215        N_means = []
216        N_stds = []
217        for N in N_values:
218            scores = np.array(random_results[N])
219            N_means.append(np.mean(scores))
220            N_stds.append(np.std(scores))
221        plt.figure(figsize=(8, 6))
222        plt.errorbar(list(N_values), N_means, yerr=N_stds, fmt='o-', capsize=5)
223        plt.xlabel('Number of Episodes (N)')
224        plt.ylabel('Mean Discounted Return')
225        plt.title('Random Model Baseline Performance')
226        plt.grid(True)
227        plt.show()
228
229
230        total_episodes_f4 = 1500
231        eval_interval_f4 = total_episodes_f4 // 20
232        dname, quality_onPi = model_based_rl(prepFrozen4, total_episodes_f4,
                eval_interval_f4,
233                                              n_repeats=2, gamma=discount,
234                                              initial_epsilon=0.5, variant='onPi')
235        eval_points_f4 = list(range(eval_interval_f4, total_episodes_f4+1,
                eval_interval_f4))
236        plot_results(eval_points_f4, quality_onPi, "Model-Based RL onPi (FrozenLake4)"
                )
237
238
239        dname, quality_RMax = model_based_rl(prepFrozen4, total_episodes_f4,
                eval_interval_f4,
240                                              n_repeats=2, gamma=discount,
241                                              initial_epsilon=0.5, variant='RMax', M
                                                  =10, R=1)
242        plot_results(eval_points_f4, quality_RMax, "Model-Based RL RMax (FrozenLake4)"
                )
```

```
243
244    total_episodes_f8 = 50000
245    eval_interval_f8 = total_episodes_f8 // 20
246    dname, quality_onPi_f8 = model_based_rl(prepFrozen8, total_episodes_f8,
           eval_interval_f8,
247                                            n_repeats=2, gamma=discount,
248                                            initial_epsilon=0.5, variant='onPi')
249    eval_points_f8 = list(range(eval_interval_f8, total_episodes_f8+1,
           eval_interval_f8))
250    plot_results(eval_points_f8, quality_onPi_f8, "Model-Based RL onPi (
           FrozenLake8)")
251
252    dname, quality_RMax_f8 = model_based_rl(prepFrozen8, total_episodes_f8,
           eval_interval_f8,
253                                            n_repeats=2, gamma=discount,
254                                            initial_epsilon=0.5, variant='RMax', M
                                                =10, R=1)
255    plot_results(eval_points_f8, quality_RMax_f8, "Model-Based RL RMax (
           FrozenLake8)")
```

Listing 2: Model_Based_RL.py

# 4 README File

```
1  # CSCI B659: Reinforcement Learning - Assignment 2
2  **Author:** LJ Huang
3  **Semester:** Spring 2025
4
5  ## Overview
6  This repository contains the code and report for Assignment 2. The assignment
      covers two main tasks:
7  1. **Task 1:** SARSA implementation in FrozenLake4, FrozenLake8, and CartPole.
8  2. **Task 2:** Model-based RL in FrozenLake4 and FrozenLake8 environments.
9
10 ## Directory Structure
11 - **VI_PI_MPI.py**
12   Contains the implementation for Task 1.
13 - **Model_Based_RL.py**
14   Contains the implementation for Task 2.
15 - **pp2starter.py**
16   The provided startup file for setting up the FrozenLake environment (including
        rewards, number of states, and actions).
17 - **hw2_report.pdf**
18   The report containing code printouts, experimental results, plots, and
        discussion of the findings.
19 - **README.md**
20   This file.
21
22 ## Dependencies
23 - Python 3.x
24 - Gymnasium (Install with `pip install gymnasium`)
25 - NumPy (Install with `pip install numpy`)
26 - Matplotlib (Install with `pip install matplotlib`)
27
28 ## Installation and Running
29 1. **Clone or Download the Repository:**
30    Clone the repository or download the zip file and extract its contents.
31
32 2. **Run the Code:**
33    \begin{verbatim}
34    python SARSA.py         # Task 1
35    python Model_Based_RL.py  # Task 2
36    \end{verbatim}
```

Listing 3: README.file