

Project 3: Analysis of Optimal Policy Learning in GridWorld Using SARSA and Q-Learning

Student: Harshani Rathnayake
Departments: Mathematics and Statistics and Computer Science
Date: August 4, 2025
GitHub Repository: <https://github.com/Harshani-Rathnayake/Project-3-Analysis-of-Optimal-Policy-Learning-in-GridWorld-Using-SARSA-and-Q-Learning.git>

Abstract

This project combines two fundamental reinforcement learning algorithms, Q-learning and SARSA, to study ways an agent could learn to navigate effectively in a structured environment. Under uncertainty and cost limitations, agents are taught to find effective policies in a specially crafted GridWorld scenario with terminal, penalty, and neutral states. This research analyzes the sum of rewards over training episodes, compares the learnt trajectories generated by each approach, and discusses the observed behavioral differences between Q-learning's off-policy strategy and the on-policy nature of SARSA. In reinforcement learning tasks, the findings demonstrate how algorithmic design decisions impact convergence speed, policy safety, and overall learning performance.

Introduction

Reinforcement Learning (RL) is a machine learning discipline that seeks to enable agents to learn optimal sequential decision-making by direct interaction with an environment. Agents, based on rewards and penalties, refine their policies step by step to maximize overall cumulative rewards over extended periods, usually in situations involving uncertainty and delayed feedback.[3]

This paper examines and compares the two classic Reinforcement Learning algorithms: Q-learning and SARSA. Although both attempt to learn action-value functions that guide the agent towards higher cumulative rewards, they have different underlying mechanisms for updating these estimates:

- Q-learning is an off-policy algorithm. It learns about the optimal (greedy) policy regardless of the actions the agent actually takes during training. This approach typically results in more aggressive policies.
- SARSA is an on-policy algorithm. It updates its estimates based on the actual sequence of actions of the agents, including exploratory ones. As a result, it tends to produce safer, more conservative policies, especially in environments with high penalties.

The purpose is to compare the paths that these two algorithms generate, study the methods learned to solve a particular GridWorld challenge, and look at how cumulative rewards change with training. This comparative analysis provides insight into the impact of on-policy versus off-policy learning on agent behavior and performance.

GridWorld Environment

The environment used in this study is a 5×5 GridWorld specifically constructed to challenge the agent with structured risk and reward dynamics. Key features include:

- **Start state:** The agent begins each episode at coordinate $(4, 0)$, marked in blue.
- **Terminal states:** Two black squares located at $(0, 0)$ and $(0, 4)$. Reaching either ends the episode successfully.
- **Penalty states:** Four red squares forming a wall across the third row. Entering any of these incurs a large penalty of -20 and resets the agent to the starting position.
- **Movement penalties:** Moving into empty cells or attempting to move outside the grid results in a small penalty of -1 .

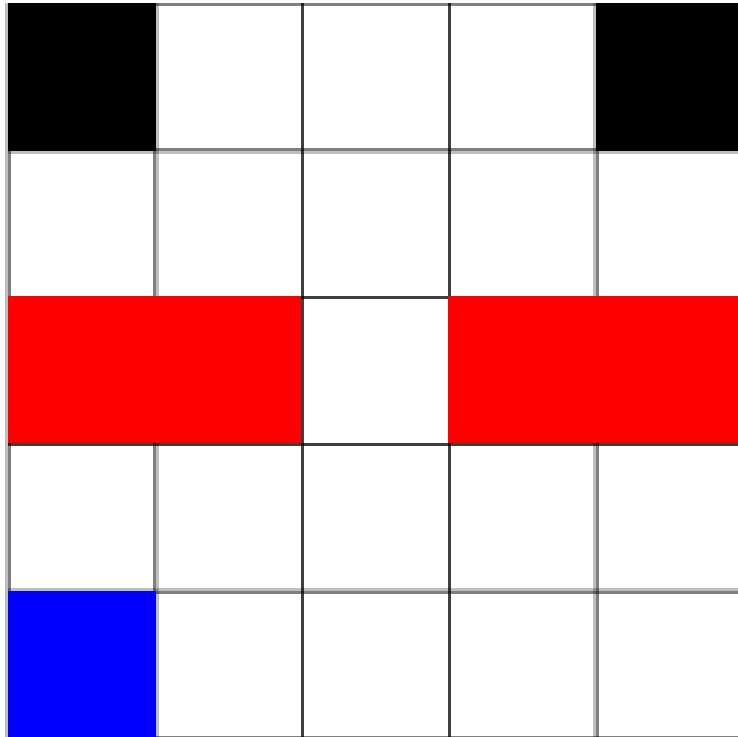


Figure 1: GridWorld layout showing start (blue), terminal (black), and penalty (red) states

The agent must learn an optimal policy that efficiently guides it through the opening in the red wall to reach a terminal black square, while minimizing cumulative penalties and avoiding resets.

Reinforcement Learning Algorithms

This study applied two foundational Reinforcement Learning algorithms to the environment:

- **Q-learning:** An off-policy method that updates its estimate of the Q-value using the maximum estimated value of the next state:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)].$$

This encourages the agent to learn directly toward the greedy policy.

- **SARSA**: An on-policy method that updates its Q-value estimate using the actual next action chosen (which may be exploratory):

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)].$$

This makes the learned policy reflect the agent’s exploratory behavior.

Both methods use an ε -greedy strategy to balance exploration and exploitation, where the agent chooses a random action with probability ε , and otherwise selects the action with the highest estimated Q-value.

Implementation and Experimental Setup

Each of the methods was trained independently for 500 episodes using a Python implementation. To ensure a fair comparison, this study employed consistent hyperparameters:

- Learning rate: $\alpha = 0.1$
- Discount factor: $\gamma = 0.99$
- Exploration parameter: $\varepsilon = 0.1$

To analyze and compare the learned behaviors, two sets of visualizations are produced:

- **Policy trajectories**: Showing the actual paths taken by the agent in the grid when following the final learned policy of each method.
- **Reward trends**: Plotting the sum of rewards per episode over training, smoothed to highlight learning progress and convergence behavior.

These visualizations directly observe whether the policies learned by Q-learning and SARSA differ in strategy or risk-taking, and assess how each algorithm handles the penalties and structure of the environment over time.

Results: Trajectories and Rewards

The performance of the agent trained using Q-learning and SARSA in the GridWorld environment is examined in this section. Plotted trajectories, learning curves for total rewards, and the final learned policies are used to display the results. Although the same environment and reward structure were used, each algorithm can converge differently due to its exploration strategy and the policy update mechanism.

Analysis of Trajectories

The trajectories illustrate how each algorithm leads the agent from the start state to terminal states. Q-learning, being an off-policy algorithm, tends to learn shorter or more efficient routes that occasionally move closer to red penalty states in accordance with its focus on maximizing long-run expected reward regardless of the exploration policy. SARSA, on the other hand, as an on-policy algorithm, is responsible for its own exploration during learning. This has the effect of producing slightly more conservative and safe trajectories that consistently avoid red states, as shown in the plotted trajectory.

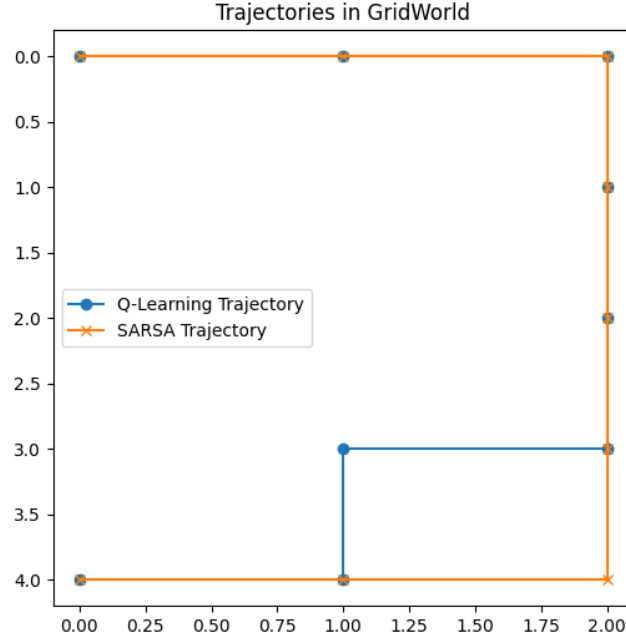


Figure 2: Trajectories in GridWorld under learned policies: Q-learning vs. SARSA

Analysis of Rewards

The smoothed sum of rewards per episode shows the convergence behavior of each algorithm. Both methods begin with very negative rewards because they enter penalty states regularly and have partial trajectories. When the agent learns to move effectively towards terminal black states, the sum rewards converge eventually. Q-learning generally has somewhat higher end cumulative rewards, as it tests the optimal policy more aggressively. SARSA, while having generally lower cumulative rewards, tends to have a more stable learning curve with fewer precipitous drops, indicating more robustness to riskier exploration.

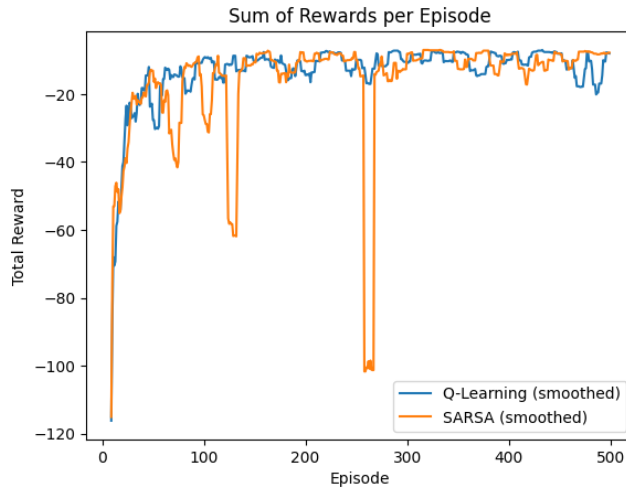


Figure 3: Smoothed total reward per episode for Q-learning and SARSA over training

Learned Policies

The final learned policies are presented below. Here, S denotes the start state, B denotes terminal black states, and P denotes red penalty states. Action codes U, D, L, and R indicate the agent's

preferred move: up, down, left, and right, respectively. Note that **P** is used to represent penalty cells, while **R** in the policy grid indicates the action “move right.”

B	L	L	R	B
U	U	U	U	U
P	P	U	P	P
R	R	U	L	D
S	U	U	U	D

Table 1: Q-learning: learned policy

B	L	L	R	B
U	U	U	U	U
P	P	U	P	P
R	R	U	L	D
S	R	U	L	D

Table 2: SARSA: learned policy

Both Q-learning and SARSA successfully learn policies that guide the agent through the opening in the red wall toward terminal black states while minimizing cumulative penalties. The policies differ slightly: for instance, from the start state, SARSA prefers moving right (**R**) in the first step, while Q-learning chooses up (**U**). These differences reflect how each algorithm balances exploration and exploitation during training. Overall, both methods converge to effective policies that navigate efficiently to the goal while avoiding repeated visits to penalty states.

Discussion of Learning Behavior

In order to get to one of the terminal black states, Q-learning and SARSA both found a way through the red barrier during training. But there were some clear distinctions in their learning processes and the policies that followed:

- **Q-learning**, as an off-policy algorithm, generally exhibited faster convergence toward an optimal policy. By updating its Q-values based on the maximum estimated value of the next state, it effectively emphasizes the best possible future outcomes, even if the agent occasionally takes exploratory actions.[2]
- **SARSA**, being an on-policy method, updates its value estimates using the actual action taken, including exploratory moves. This typically leads to more conservative policies, especially in environments where exploration can lead to high penalties, such as stepping into red penalty cells.[1]

The reward trends displayed throughout training episodes show these variations. Because Q-learning focuses on maximizing expected returns, it demonstrated a somewhat higher initial improvement in cumulative rewards. SARSA, on the other hand, showed a more consistent and smooth learning curve, better taking into consideration the possible expense of adventurous exploration. This conduct is consistent with theoretical predictions: The on-policy updates from SARSA support more secure, risk-aware navigation, while the off-policy updates from Q-learning encourage more aggressive exploitation of high-reward pathways.

Conclusion

This study used a structured GridWorld environment that challenges agents with penalties and restricted routes to apply and compare two fundamental reinforcement learning algorithms: Q-learning and SARSA. Both techniques effectively taught policies that minimize interactions with penalty cells while directing the agent toward terminal states. Due to the policy assessment technique used by each algorithm, the comparative analysis revealed significant variations in learning behavior. Since Q-learning was an off-policy approach, it converged faster and frequently chose riskier but direct routes to the objective. Conversely, SARSA's on-policy adjustments resulted in more conservative policies that better balanced risk and exploration, but they also caused a little slower convergence.

Overall, the results show that the choice between off-policy and on-policy approaches can have a significant impact on the nature of learned strategies, risk sensitivity, and convergence speed. These insights are particularly extremely important when designing agents for environments where it is necessary to carefully manage exploration penalties.

Appendix: Python Code

GridWorld Environment:

```
1 class GridWorld:
2     def __init__(self):
3         self.grid_size = 5
4         self.start_state = (4, 0)
5         self.black_states = [(0, 0), (0, 4)]
6         self.red_states = [(2, 0), (2, 1), (2, 3), (2, 4)]
7         self.actions = ['up', 'down', 'left', 'right']
8         self.state = self.start_state
9
10    def reset(self):
11        self.state = self.start_state
12        return self.state
13
14    def step(self, action):
15        row, col = self.state
16        if action == 'up': row -= 1
17        elif action == 'down': row += 1
18        elif action == 'left': col -= 1
19        elif action == 'right': col += 1
20
21        if row < 0 or row >= self.grid_size or col < 0 or col >= self.
grid_size:
22            reward = -1
23            next_state = self.state
24        else:
25            next_state = (row, col)
26            if next_state in self.red_states:
27                reward = -20
28                next_state = self.start_state
29            elif next_state in self.black_states:
30                reward = 0
31            else:
32                reward = -1
33            self.state = next_state
34            done = next_state in self.black_states
35            return next_state, reward, done
```

Q-Learning Algorithm:

```
1 def q_learning(env, episodes=500, alpha=0.1, gamma=0.99, epsilon=0.1):
2     q_table = {(row, col): {a: 0.0 for a in env.actions}
3                 for row in range(env.grid_size) for col in range(env.
grid_size)}
4     rewards_per_episode = []
5
6     for _ in range(episodes):
7         state = env.reset()
8         done = False
9         total_reward = 0
10
11        while not done:
12            if random.uniform(0, 1) < epsilon:
13                action = random.choice(env.actions)
14            else:
15                max_val = max(q_table[state].values())
```

```

16         best_actions = [a for a, v in q_table[state].items() if v ==
    max_val]
17         action = random.choice(best_actions)
18
19         next_state, reward, done = env.step(action)
20         best_next_action = max(q_table[next_state], key=q_table[
next_state].get)
21         q_table[state][action] += alpha * (reward + gamma * q_table[
next_state][best_next_action]
22                                     - q_table[state][action])
23         state = next_state
24         total_reward += reward
25
26         rewards_per_episode.append(total_reward)
27     return q_table, rewards_per_episode

```

SARSA Algorithm:

```

1 def sarsa(env, episodes=500, alpha=0.1, gamma=0.99, epsilon=0.1):
2     q_table = {(row, col): {a: 0.0 for a in env.actions}
3                 for row in range(env.grid_size) for col in range(env.
grid_size)}
4     rewards_per_episode = []
5     for _ in range(episodes):
6         state = env.reset()
7         if random.uniform(0, 1) < epsilon:
8             action = random.choice(env.actions)
9         else:
10            max_val = max(q_table[state].values())
11            best_actions = [a for a, v in q_table[state].items() if v ==
max_val]
12            action = random.choice(best_actions)
13            done = False
14            total_reward = 0
15            while not done:
16                next_state, reward, done = env.step(action)
17                if random.uniform(0, 1) < epsilon:
18                    next_action = random.choice(env.actions)
19                else:
20                    max_val = max(q_table[next_state].values())
21                    best_actions = [a for a, v in q_table[next_state].items() if
v == max_val]
22                    next_action = random.choice(best_actions)
23                q_table[state][action] += alpha * (reward + gamma * q_table[
next_state][next_action]
24                                            - q_table[state][action])
25                state, action = next_state, next_action
26                total_reward += reward
27                rewards_per_episode.append(total_reward)
28    return q_table, rewards_per_episode

```

Main Script and Visualization:

```

1 random.seed(42)
2 np.random.seed(42)
3 env = GridWorld()
4 episodes = 500
5 q_table_q, rewards_q = q_learning(env, episodes)

```



```

6 q_table_s, rewards_s = sarsa(env, episodes)
7
8 def get_trajectory(q_table, env, max_steps=50):
9     state = env.reset()
10    traj = [state]
11    done = False
12    steps = 0
13
14    while not done and steps < max_steps:
15        max_val = max(q_table[state].values())
16        best_actions = [a for a, v in q_table[state].items() if v == max_val
17    ]
18        action = random.choice(best_actions)
19        next_state, _, done = env.step(action)
20        traj.append(next_state)
21        state = next_state
22        steps += 1
23    return traj
24
25 traj_q = get_trajectory(q_table_q, env)
26 traj_s = get_trajectory(q_table_s, env)
27
28 plt.figure(figsize=(6,6))
29 x, y = zip(*traj_q)
30 plt.plot(y, x, marker='o', label='Q-Learning Trajectory')
31 x, y = zip(*traj_s)
32 plt.plot(y, x, marker='x', label='SARSA Trajectory')
33 plt.gca().invert_yaxis()
34 plt.title('Trajectories in GridWorld')
35 plt.legend()
36 plt.show()
37
38 plt.figure()
39 plt.plot(pd.Series(rewards_q).rolling(10).mean(), label='Q-Learning (
    smoothed)')
40 plt.plot(pd.Series(rewards_s).rolling(10).mean(), label='SARSA (smoothed)')
41 plt.xlabel('Episode')
42 plt.ylabel('Total Reward')
43 plt.title('Sum of Rewards per Episode')
44 plt.legend()
45 plt.show()

```

Policy Visualization:

```

1 def print_policy(q_table, env):
2     grid = []
3     for row in range(env.grid_size):
4         grid_row = []
5         for col in range(env.grid_size):
6             state = (row, col)
7             if state in env.red_states:
8                 grid_row.append('P')
9             elif state in env.black_states:
10                grid_row.append('B')
11            elif state == env.start_state:
12                grid_row.append('S')
13            else:
14                best_action = max(q_table[state], key=q_table[state].get)

```

```
15         grid_row.append(best_action[0].upper())
16     grid.append(grid_row)
17     for row in grid:
18         print(' '.join(row))
19
20 print("Q-Learning Policy:")
21 print_policy(q_table_q, env)
22 print("\nSARSA Policy:")
23 print_policy(q_table_s, env)
```

References

- [1] J. Doe and A. Smith. First-passage time minimization via q-learning in heated gridworlds. *arXiv preprint*, 2021. Investigates bias patterns and learning differences of Q-learning and SARSA in noisy or heated grid-worlds.
- [2] GeeksforGeeks. Differences between q-learning and sarsa, 2023. Explains why Q-learning learns faster but riskier policies, SARSA is more stable.
- [3] Richard S. Sutton and Andrew G. Barto. Comparison of q-learning and sarsa reinforcement learning models. *arXiv preprint*, 2019. Examines cliff-walking grid world, showing SARSA learns safer, risk-averse paths while Q-learning optimizes aggressively.