

```
In [1]: 1 # Install relevant Libraries
2 !pip install numpy matplotlib tqdm scipy
```

```
Requirement already satisfied: numpy in c:\users\student\.conda\envs\dsai\lib\site-packages (1.24.3)
Requirement already satisfied: matplotlib in c:\users\student\.conda\envs\dsai\lib\site-packages (3.9.2)
Requirement already satisfied: tqdm in c:\users\student\.conda\envs\dsai\lib\site-packages (4.66.5)
Requirement already satisfied: scipy in c:\users\student\.conda\envs\dsai\lib\site-packages (1.14.1)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: cycler>=0.10 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: kiwisolver>=1.3.1 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: packaging>=20.0 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: pillow>=8 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib) (10.
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\student\.conda\envs\dsai\lib\site-packages (from matplotlib)
Requirement already satisfied: colorama in c:\users\student\.conda\envs\dsai\lib\site-packages (from tqdm) (0.4.6)
Requirement already satisfied: six>=1.5 in c:\users\student\.conda\envs\dsai\lib\site-packages (from python-dateutil>=
2.7->matplotlib) (1.16.0)
```

```
In [2]: 1
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from tqdm import tqdm
5 from IPython.display import clear_output
6 %matplotlib inline
```

Problem Statement

In this section we will implement tabular SARSA and Q-learning algorithms for a grid world navigation task.

Environment details

The agent can move from one grid coordinate to one of its adjacent grids using one of the four actions: UP, DOWN, LEFT and RIGHT. The goal is to go from a randomly assigned starting position to goal position.

Actions that can result in taking the agent off the grid will not yield any effect. Lets look at the environment.

```
In [3]: 1 DOWN = 0
2 UP = 1
3 LEFT = 2
4 RIGHT = 3
5 actions = [DOWN, UP, LEFT, RIGHT]
```

Let us construct a grid in a text file.

```
In [4]: 1 !type grid_world2.txt
```

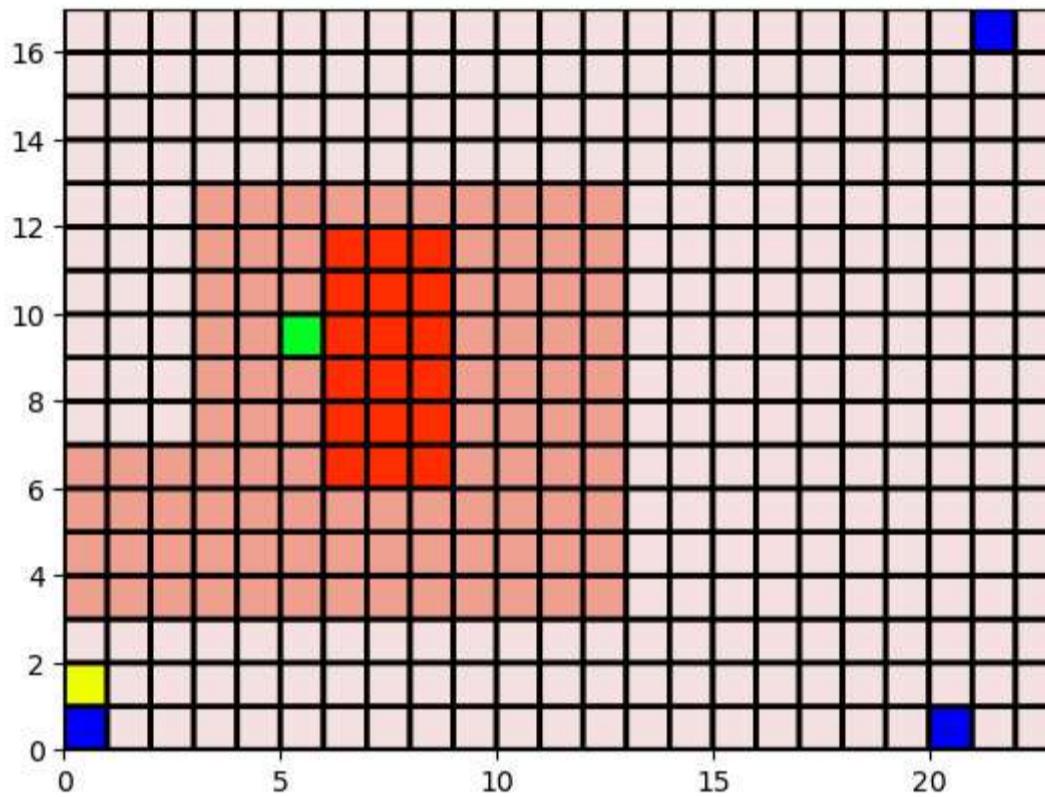
```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 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 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
0 0 0 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
```

This is a 17×23 grid. The reward when an agent goes to a cell is negative of the value in that position in the text file (except if it is the goal cell). We will define the goal reward as 100. We will also fix the maximum episode length to 10000.

Now let's make it more difficult. We add stochasticity to the environment: with probability 0.2 agent takes a random action (which can be other than the chosen action). There is also a westerly wind blowing (to the right). Hence, after every time-step, with probability 0.5 the agent also moves an extra step to the right.

Now let's plot the grid world.

```
In [8]: 1 world = 'grid_world2.txt'
2 goal_reward = 100
3 start_states = [(0,0), (0,20), (16,21)]
4 goal_states=[(9,5)]
5 max_steps=10000
6
7 from grid_world import GridWorldEnv, GridWorldWindyEnv
8
9 env = GridWorldEnv(world, goal_reward=goal_reward, start_states=start_states, goal_states=goal_states,
10                      max_steps=max_steps, action_fail_prob=0.2)
11 plt.figure(figsize=(10, 10))
12 plt.show()
13 # Go UP
14 env.step(UP)
15 env.render(ax=plt, render_agent=True)
```



<Figure size 1000x1000 with 0 Axes>

Legend

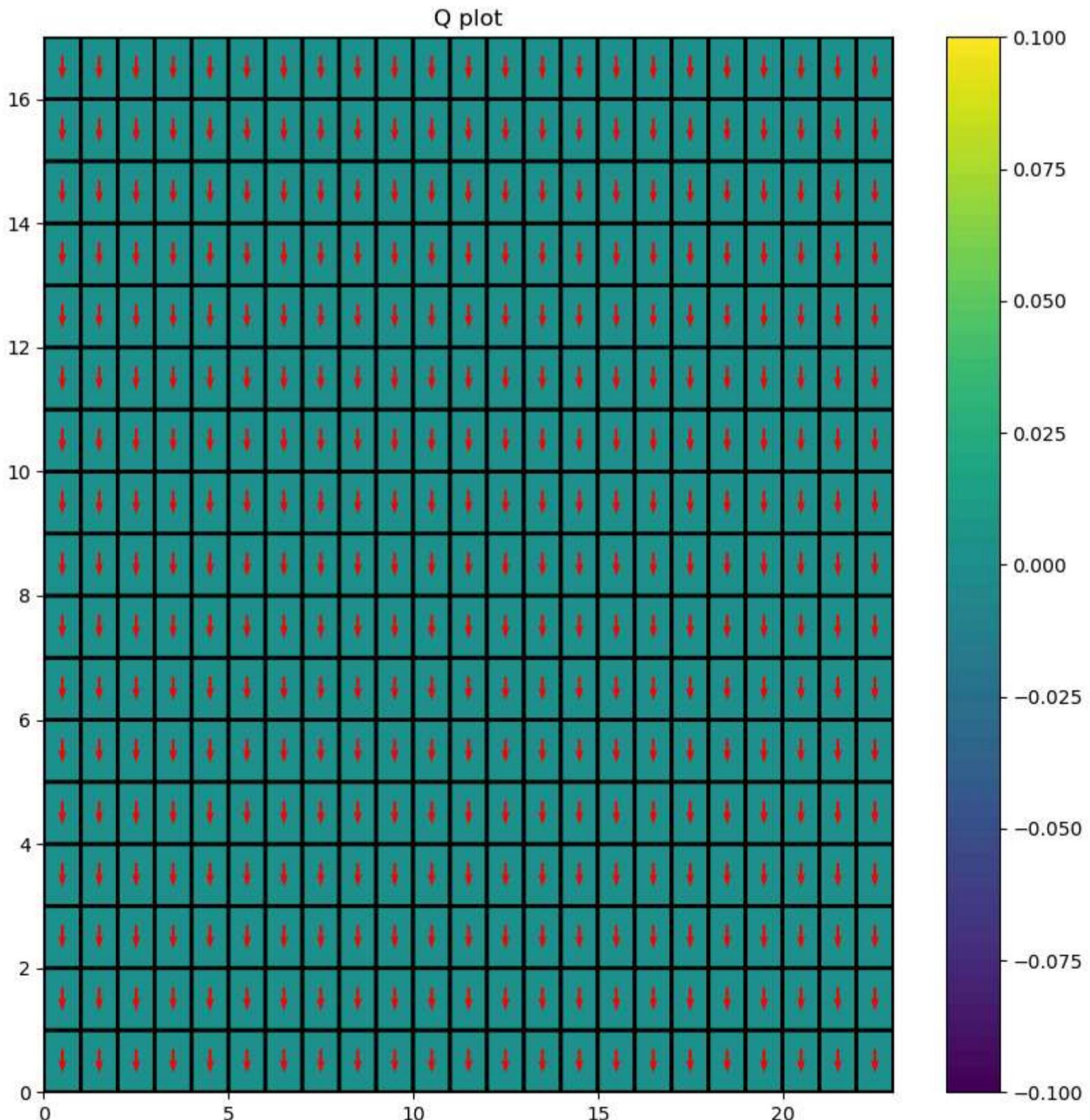
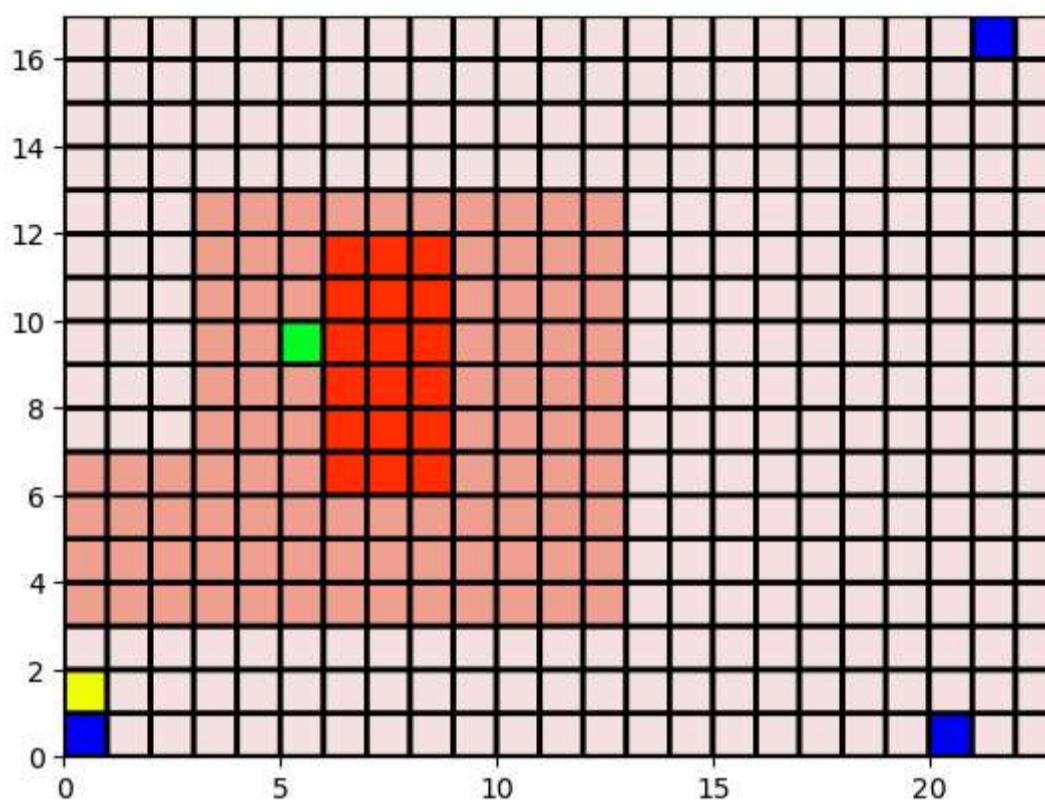
- **Blue** is the **start state**.
- **Green** is the **goal state**.
- **Yellow** is current **state of the agent**.
- **Redness** denotes the extent of **negative reward**.

Q values

We can use a 3D array to represent Q values. The first two indices are X, Y coordinates and last index is the action.

In [9]:

```
1 from grid_world import plot_Q
2
3 Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
4
5 plot_Q(Q)
6
7 Q.shape
```



Out[9]: (17, 23, 4)

Exploration strategies

1. Epsilon-greedy
2. Softmax

In [10]:

```

1 from scipy.special import softmax
2
3 seed = 42
4 rg = np.random.RandomState(seed)
5
6 # Epsilon greedy
7 def choose_action_epsilon(Q, state, epsilon, rg=rg):
8     if not Q[state[0], state[1]].any() or rg.rand() < epsilon:
9         return rg.choice(Q.shape[-1])
10    else:
11        return np.argmax(Q[state[0], state[1]])
12
13 # Softmax
14 def choose_action_softmax(Q, state, rg=rg):
15     return rg.choice(Q.shape[-1], p = softmax(Q[state[0]], state[1]))

```

SARSA

Now we implement the SARSA algorithm.

Recall the update rule for SARSA:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Hyperparameters

So we have some hyperparameters for the algorithm:

- α
- number of episodes.
- ϵ : For epsilon greedy exploration

In [11]:

```

1 # initialize Q-value
2 Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
3
4 alpha0 = 0.4
5 gamma = 0.9
6 episodes = 10000
7 epsilon0 = 0.1

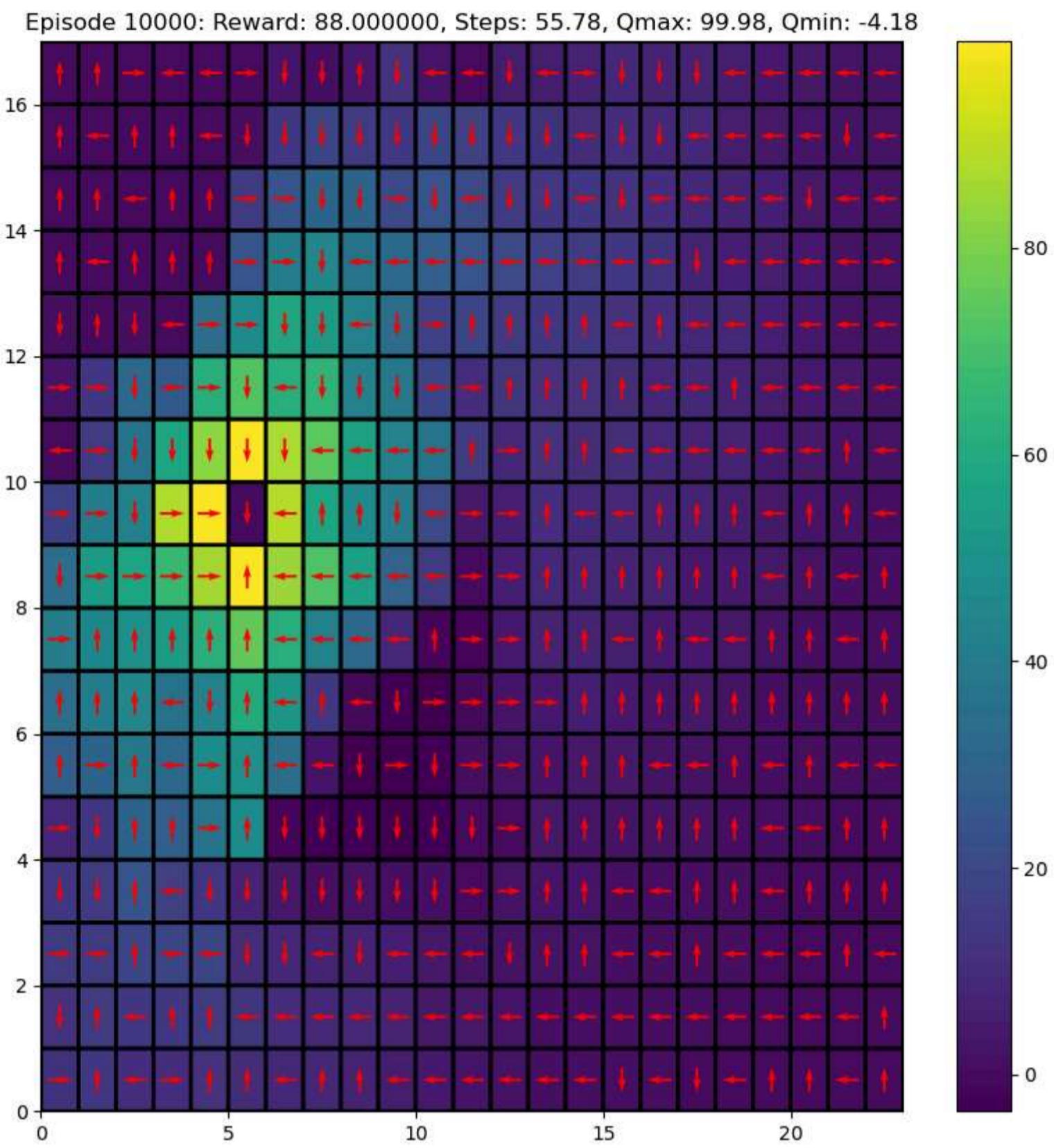
```

Let's implement SARSA

In [12]:

```
1 print_freq = 10
2
3 def sarsa(env, Q, gamma = 0.9, plot_heat = False, choose_action = choose_action_softmax):
4
5     episode_rewards = np.zeros(episodes)
6     steps_to_completion = np.zeros(episodes)
7     if plot_heat:
8         clear_output(wait=True)
9         plot_Q(Q)
10    epsilon = epsilon0
11    alpha = alpha0
12    for ep in tqdm(range(episodes)):
13        tot_reward, steps = 0, 0
14
15        # Reset environment
16        state = env.reset()
17        action = choose_action(Q, state)
18        done = False
19        while not done:
20            state_next, reward, done = env.step(action)
21            action_next = choose_action(Q, state_next)
22
23            # update equation
24            Q[state[0], state[1], action] += alpha*(reward + gamma*Q[state_next[0], state_next[1], action_next] -
25                                            Q[state[0], state[1], action])
26
27            tot_reward += reward
28            steps += 1
29
30            state, action = state_next, action_next
31
32        episode_rewards[ep] = tot_reward
33        steps_to_completion[ep] = steps
34
35        if (ep+1)%print_freq == 0 and plot_heat:
36            clear_output(wait=True)
37            plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f, Qmax: %.2f, Qmin: %.2f" %(
38                ep+1, np.mean(episode_rewards[ep-print_freq+1:ep]),
39                np.mean(steps_to_completion[ep-print_freq+1:ep]),
40                Q.max(), Q.min()))
41
42    return Q, episode_rewards, steps_to_completion
```

```
In [13]: 1 Q, rewards, steps = sarsa(env, Q, gamma = gamma, plot_heat=True, choose_action=choose_action_softmax)
```



Visualizing the policy

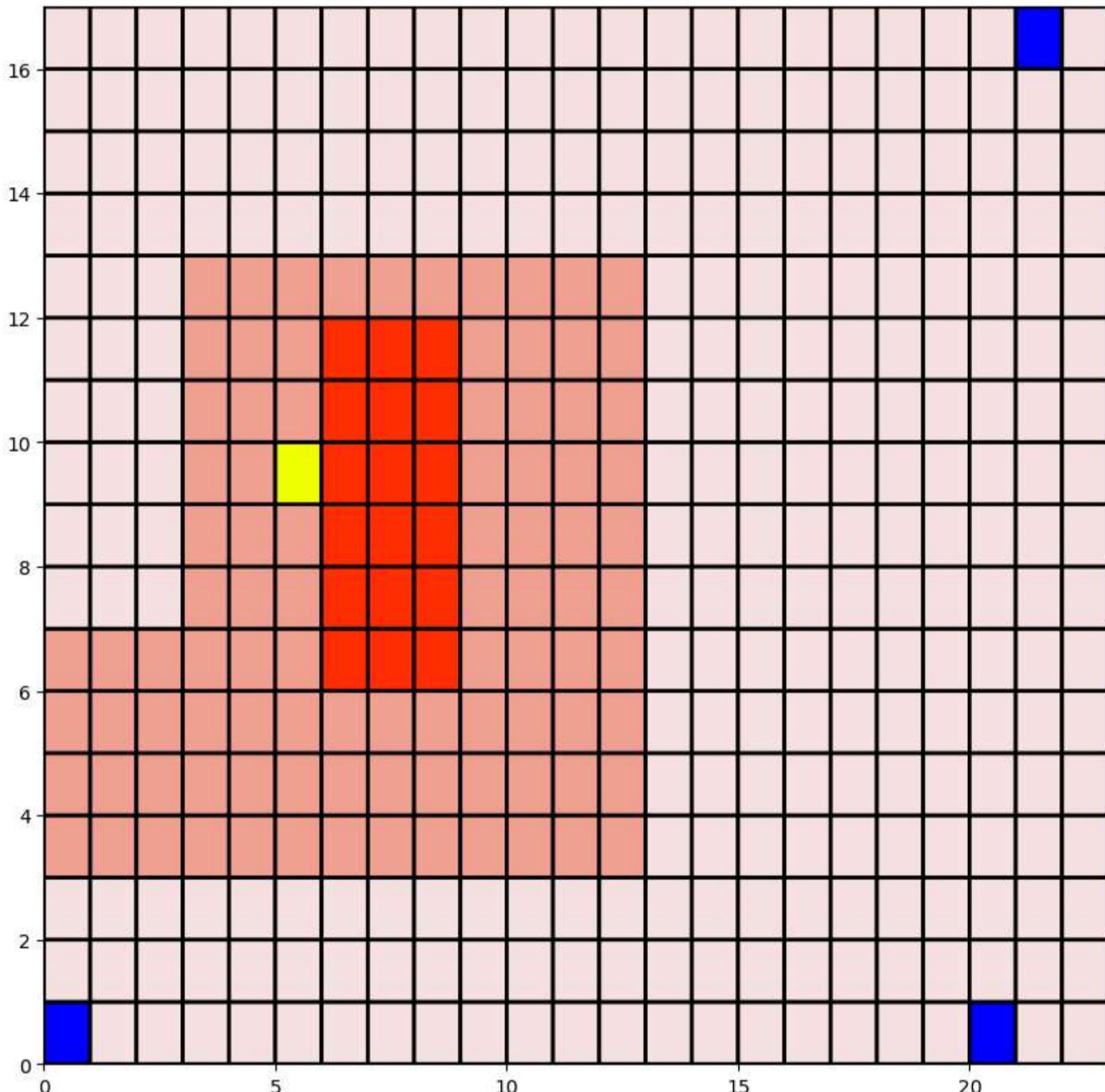
Now let's see the agent in action. Run the below cell (as many times) to render the policy;

In [14]:

```

1 from time import sleep
2
3 state = env.reset()
4 done = False
5 steps = 0
6 tot_reward = 0
7 while not done:
8     clear_output(wait=True)
9     state, reward, done = env.step(Q[state[0], state[1]].argmax())
10    plt.figure(figsize=(10, 10))
11    env.render(ax=plt, render_agent=True)
12    plt.show()
13    steps += 1
14    tot_reward += reward
15    sleep(0.2)
16 print("Steps: %d, Total Reward: %d"%(steps, tot_reward))

```



Steps: 20, Total Reward: 89

Analyzing performance of the policy

We use two metrics to analyze the policies:

1. Average steps to reach the goal
2. Total rewards from the episode

To ensure, we account for randomness in environment and algorithm (say when using epsilon-greedy exploration), we run the algorithm for multiple times and use the average of values over all runs.

In [15]:

```
1 Q_avgs, reward_avgs, steps_avgs = [], [], []
2 num_expts = 5
3
4 for i in range(num_expts):
5     print("Experiment: %d" %(i+1))
6     Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
7     rg = np.random.RandomState(i)
8     Q, rewards, steps = sarsa(env, Q)
9     Q_avgs.append(Q.copy())
10    reward_avgs.append(rewards)
11    steps_avgs.append(steps)
```

Experiment: 1

100%|██████████| 10000/10000 [00:09<00:00, 1025.89it/s]

Experiment: 2

100%|██████████| 10000/10000 [00:09<00:00, 1068.88it/s]

Experiment: 3

100%|██████████| 10000/10000 [00:14<00:00, 702.39it/s]

Experiment: 4

100%|██████████| 10000/10000 [00:10<00:00, 943.97it/s]

Experiment: 5

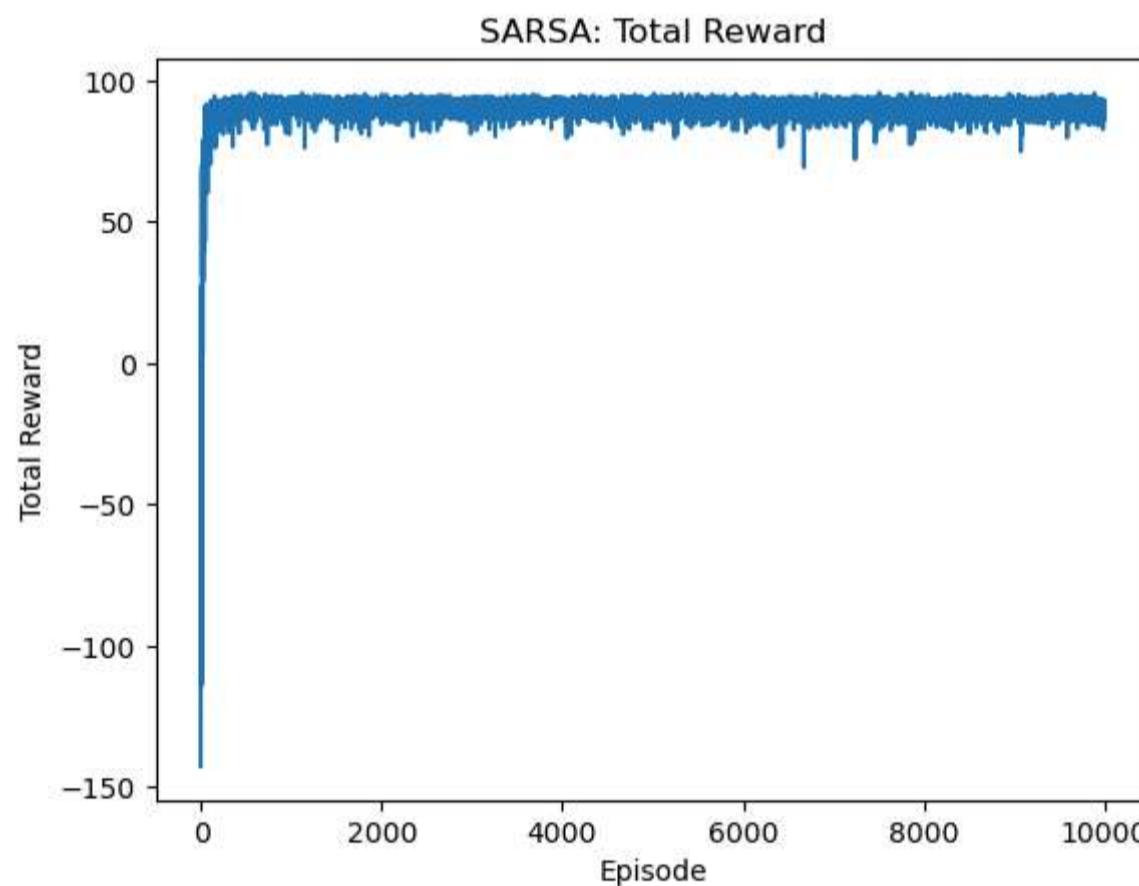
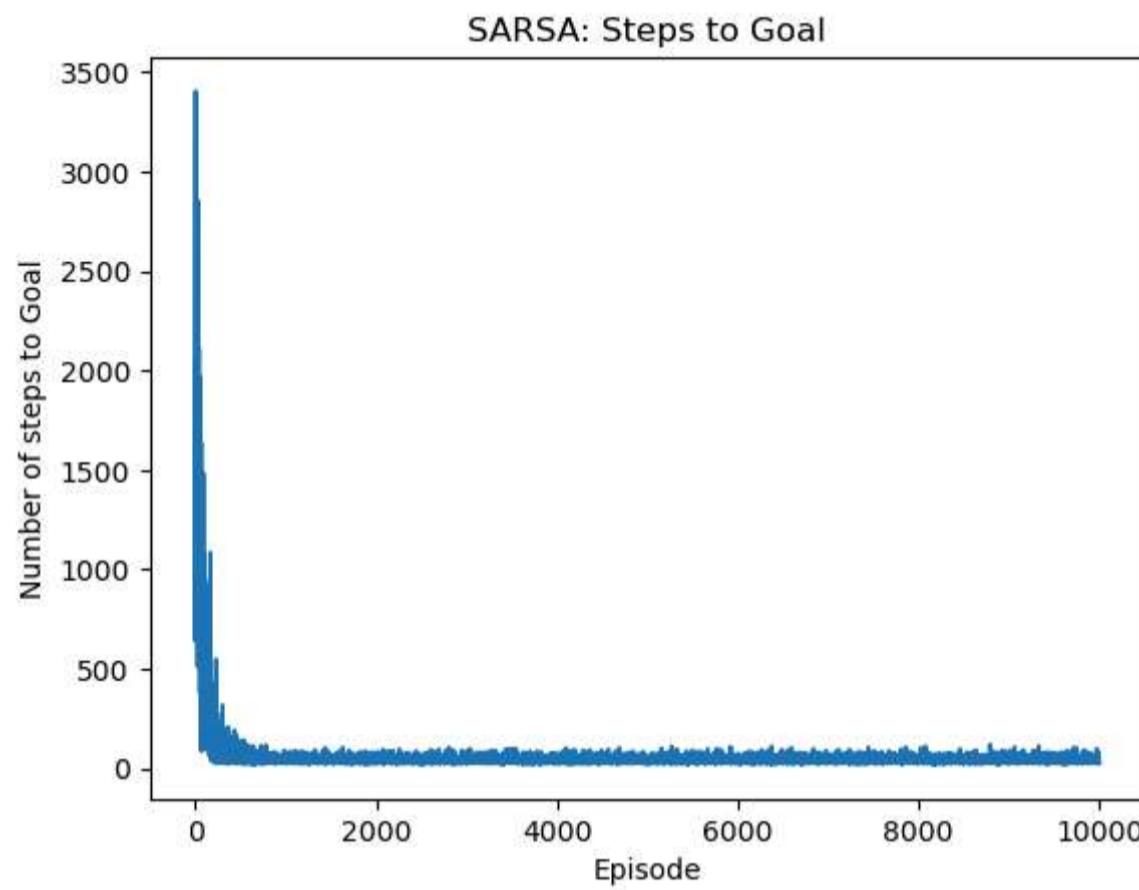
100%|██████████| 10000/10000 [00:10<00:00, 979.99it/s]

In [16]:

```

1 plt.xlabel('Episode')
2 plt.ylabel('Number of steps to Goal')
3 plt.title('SARSA: Steps to Goal')
4 plt.plot(np.arange(episodes),np.average(steps_avgs, 0))
5 plt.show()
6 plt.xlabel('Episode')
7 plt.ylabel('Total Reward')
8 plt.title('SARSA: Total Reward')
9 plt.plot(np.arange(episodes),np.average(reward_avgs, 0))
10 plt.show()

```



TODO: Q-Learning

Now, implement the Q-Learning algorithm as an exercise.

Recall the update rule for Q-Learning:

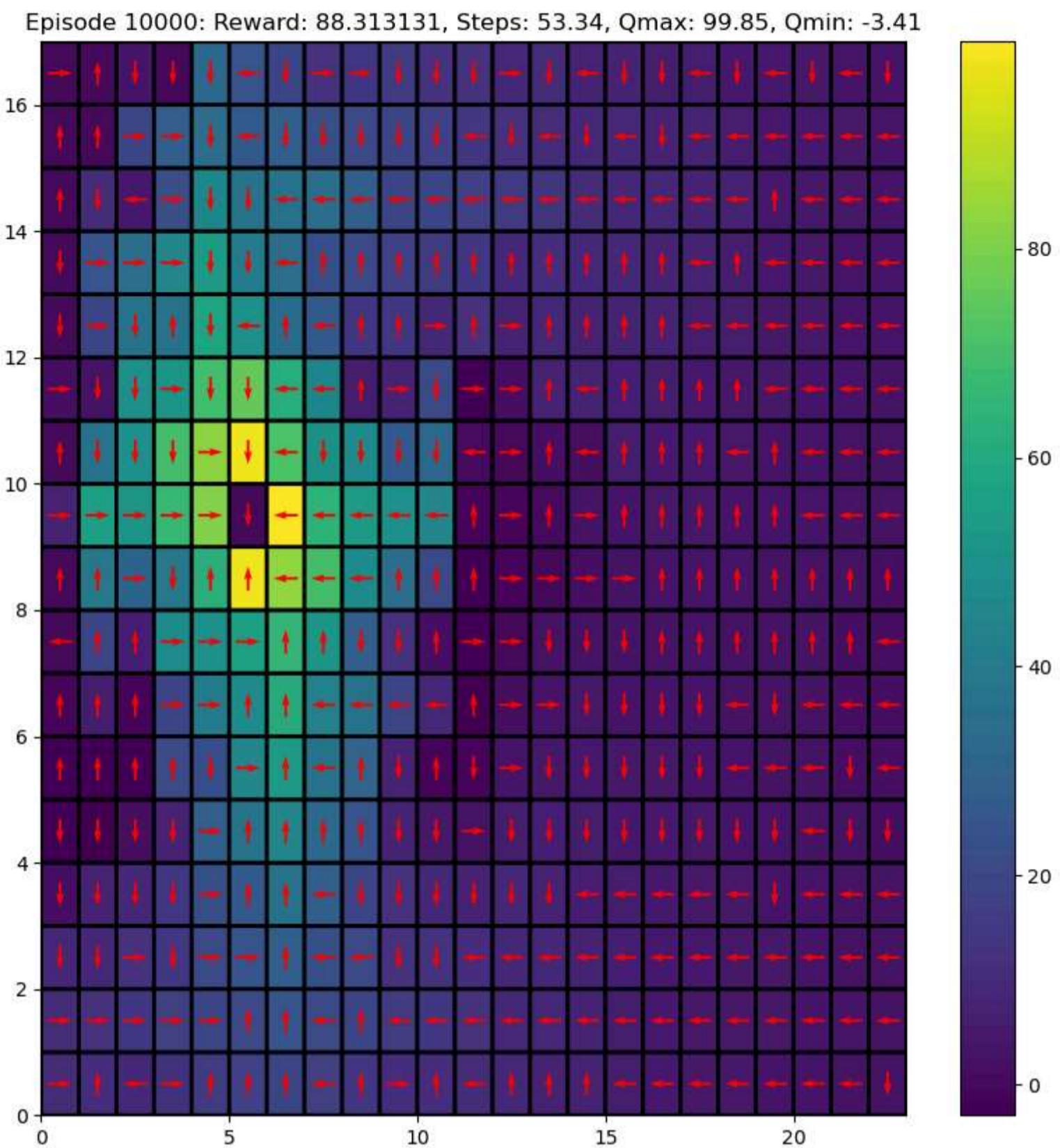
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

Visualize and compare results with SARSA. Repeat experiments 5 times (as in SARSA) and report average number of steps to goal and average episodic reward. You may use the same hyperparameters as SARSA.

In [17]:

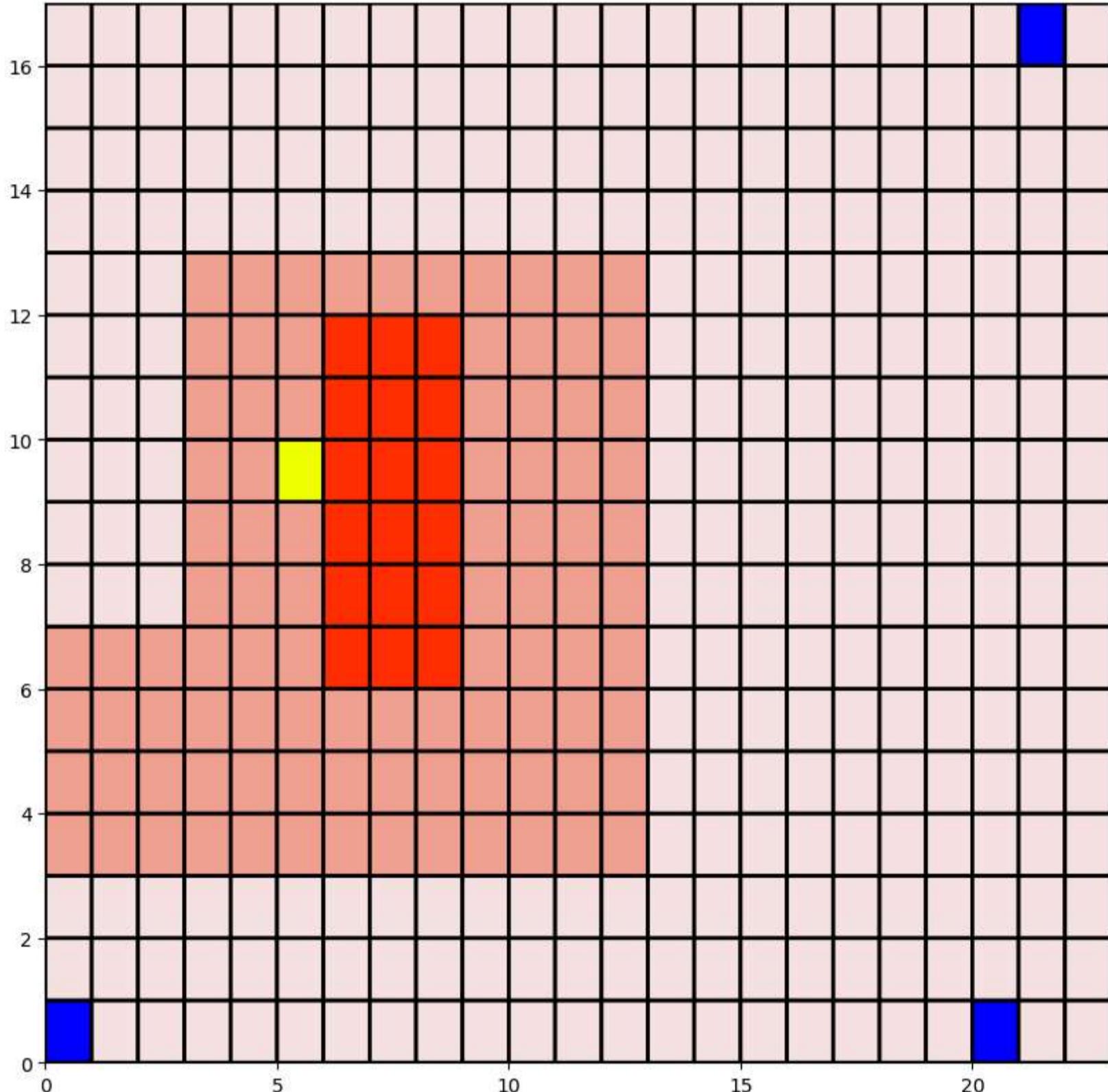
```
1 ### WRITE CODE HERE ### (Q-Learning)
2 print_freq = 100
3
4 def QLearning(env, Q, gamma = 0.9, plot_heat = False, choose_action = choose_action_softmax):
5
6     episode_rewards = np.zeros(episodes)
7     steps_to_completion = np.zeros(episodes)
8     if plot_heat:
9         clear_output(wait=True)
10        plot_Q(Q)
11
12    epsilon = epsilon0
13    alpha = alpha0
14    for ep in tqdm(range(episodes)):
15        tot_reward, steps = 0, 0
16
17        state = env.reset()
18        done = False
19        while not done:
20            action = choose_action(Q, state) # -> softmax
21            #action = choose_action(Q, state, epsilon) # -> Epsilon greedy
22            state_next, reward, done = env.step(action)
23
24            Q[state[0], state[1], action] += alpha*(reward + gamma*np.max(Q[state_next[0], state_next[1], :]) -
25                                              Q[state[0], state[1], action])
26            tot_reward += reward
27            steps += 1
28
29            state = state_next
30
31        episode_rewards[ep] = tot_reward
32        steps_to_completion[ep] = steps
33
34        if (ep+1)%print_freq == 0 and plot_heat:
35            clear_output(wait=True)
36            plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f, Qmax: %.2f, Qmin: %.2f"%
37                                         (ep+1, np.mean(episode_rewards[ep-print_freq+1:ep]),
38                                          np.mean(steps_to_completion[ep-print_freq+1:ep]),
39                                          Q.max(), Q.min()))
39
40    return Q, episode_rewards, steps_to_completion
41
42
```

```
In [18]: 1 Q, rewards, steps = QLearning(env, Q, gamma = gamma, plot_heat=True, choose_action=choose_action_softmax)
```



In [19]:

```
1 ### WRITE CODE HERE ### (visualize policy)
2 from time import sleep
3
4 state = env.reset()
5 done = False
6 steps = 0
7 tot_reward = 0
8 while not done:
9     clear_output(wait=True)
10    state, reward, done = env.step(Q[state[0], state[1]].argmax())
11    plt.figure(figsize=(10, 10))
12    env.render(ax=plt, render_agent=True)
13    plt.show()
14    steps += 1
15    tot_reward += reward
16    sleep(0.2)
17 print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
18
```



Steps: 31, Total Reward: 88

In [20]:

```

1 ### WRITE CODE HERE ### (plot metrics avged over 5 runs)
2
3 Q_avgs_Qlearning, reward_avgs_Qlearning, steps_avgs_Qlearning = [], [], []
4 num_expts = 5
5
6 for i in range(num_expts):
7     print("Experiment: %d" % (i+1))
8     Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
9     rg = np.random.RandomState(i)
10    Q, rewards, steps = sarsa(env, Q)
11    Q_avgs_Qlearning.append(Q.copy())
12    reward_avgs_Qlearning.append(rewards)
13    steps_avgs_Qlearning.append(steps)
14
15 plt.xlabel('Episode')
16 plt.ylabel('Number of steps to Goal')
17 plt.title('Q-Learning: Steps to Goal')
18 plt.plot(np.arange(episodes), np.average(steps_avgs_Qlearning, 0), color='red')
19 plt.show()
20
21 plt.xlabel('Episode')
22 plt.ylabel('Total Reward')
23 plt.title('Q-Learning: Total Reward')
24 plt.plot(np.arange(episodes), np.average(reward_avgs_Qlearning, 0), color='red')
25 plt.show()

```

Experiment: 1

100%|██████████| 10000/10000 [00:16<00:00, 599.78it/s]

Experiment: 2

100%|██████████| 10000/10000 [00:10<00:00, 916.98it/s]

Experiment: 3

100%|██████████| 10000/10000 [00:10<00:00, 920.84it/s]

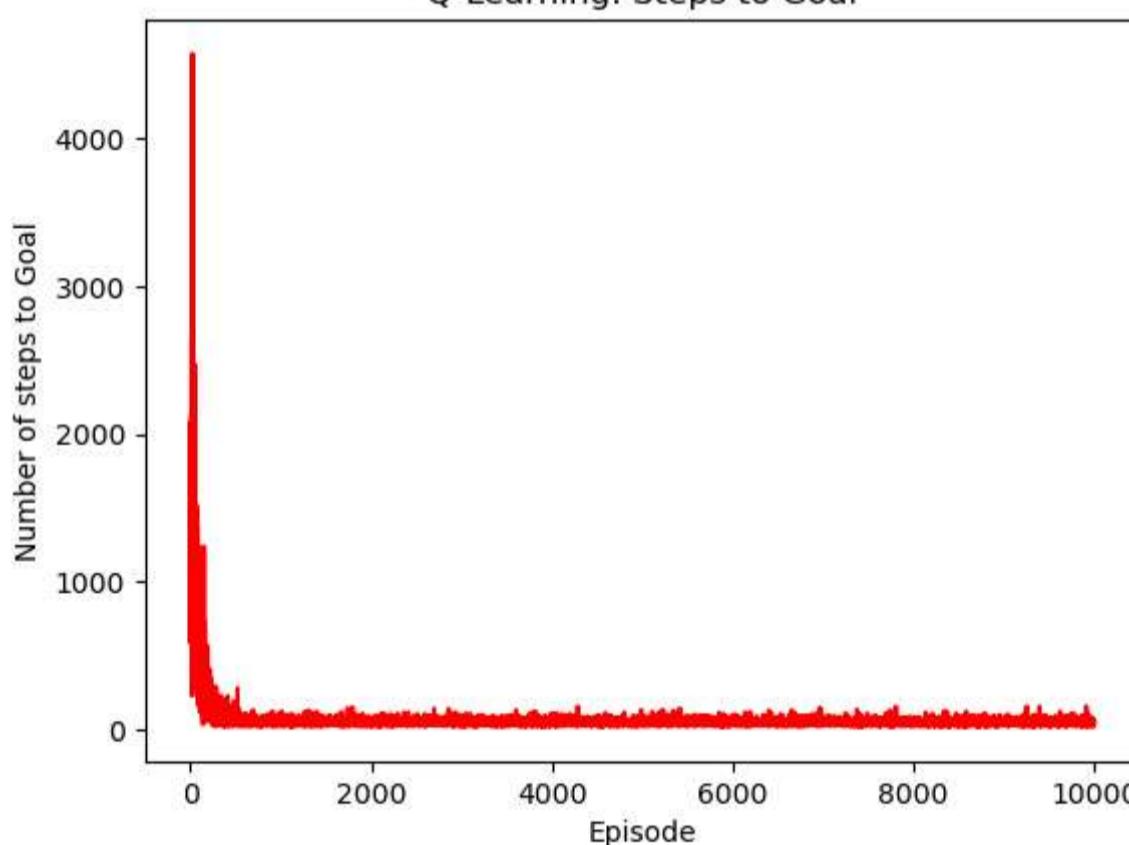
Experiment: 4

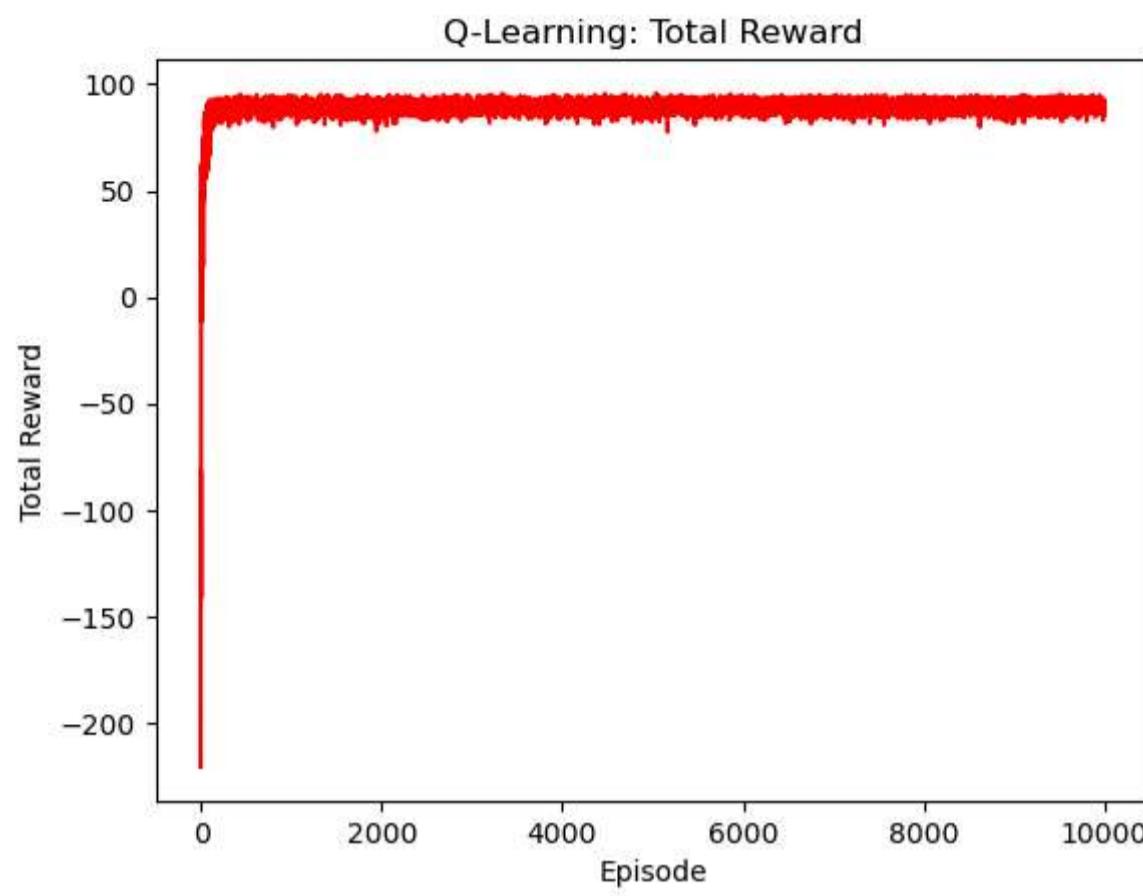
100%|██████████| 10000/10000 [00:09<00:00, 1013.53it/s]

Experiment: 5

100%|██████████| 10000/10000 [00:13<00:00, 763.99it/s]

Q-Learning: Steps to Goal





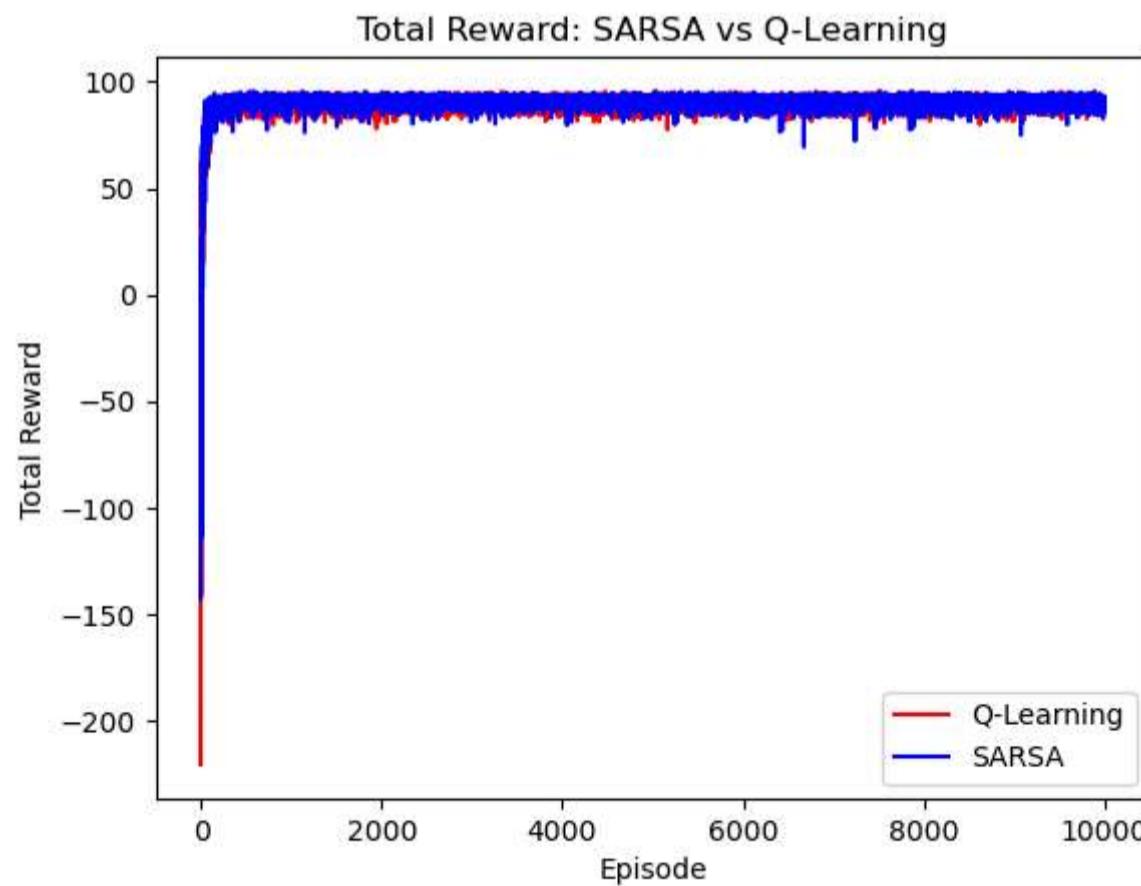
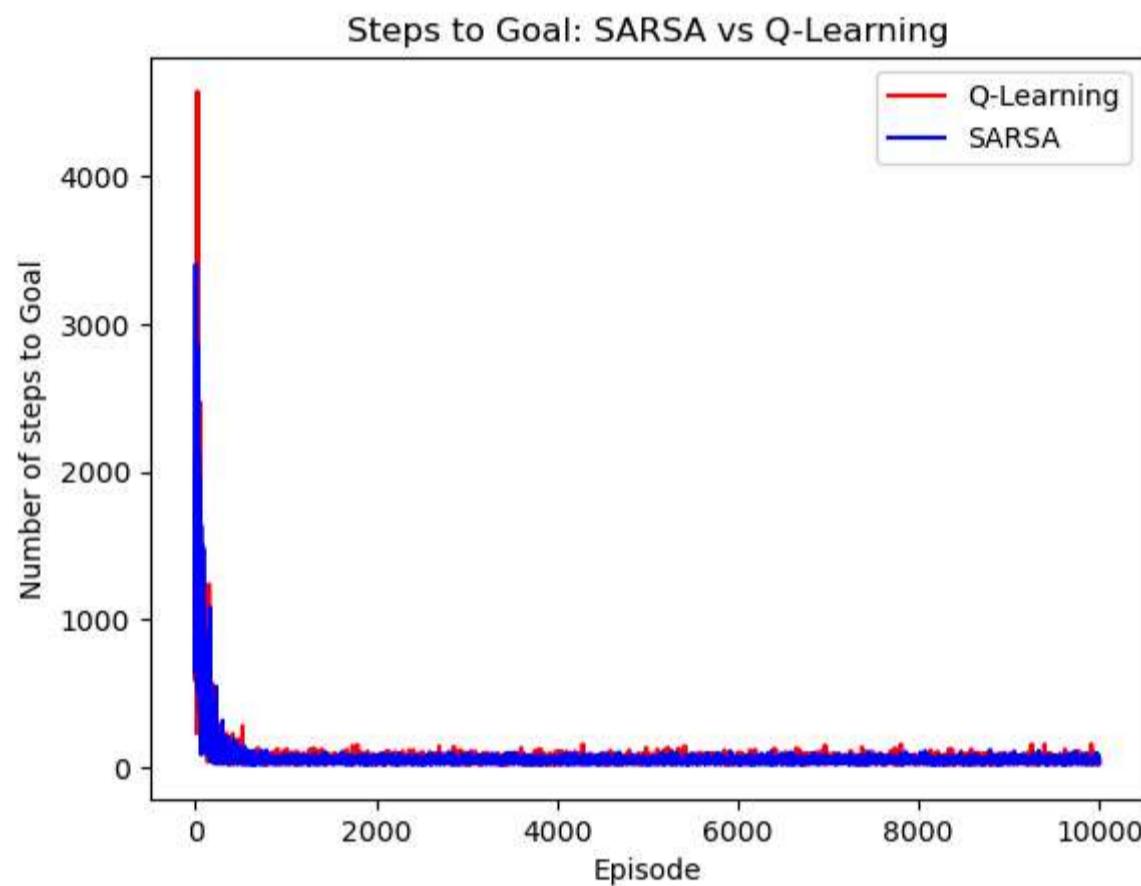
TODO: Compare the policies learnt by Q Learning and SARSA and note differences (if any)

In [21]:

```

1 ### WRITE CODE HERE ### (compare Q-Learning and SARSA performance)
2
3 plt.xlabel('Episode')
4 plt.ylabel('Number of steps to Goal')
5 plt.title('Steps to Goal: SARSA vs Q-Learning')
6 plt.plot(np.arange(episodes), np.average(steps_avgs_Qlearning, 0), color='red', label='Q-Learning')
7 plt.plot(np.arange(episodes), np.average(steps_avgs, 0), color='blue', label='SARSA')
8 plt.legend()
9 plt.show()
10
11
12 plt.xlabel('Episode')
13 plt.ylabel('Total Reward')
14 plt.title('Total Reward: SARSA vs Q-Learning')
15 plt.plot(np.arange(episodes), np.average(reward_avgs_Qlearning, 0), color='red', label='Q-Learning')
16 plt.plot(np.arange(episodes), np.average(reward_avgs, 0), color='blue', label='SARSA')
17 plt.legend()
18 plt.show()
19

```



Observations:

1. SARSA and Qlearning algorithms achieve similar long term rewards indicating that the algorithms converge to near optimal polices.
2. Q-Learning, tends to take slightly more steps on average as it aggressively updates based on the maximum Q values, leading to riskier exploration. In contrast, SARSA's on policy update results in a more conservative approach, often leading to convergence with fewer steps as observed in the plot.

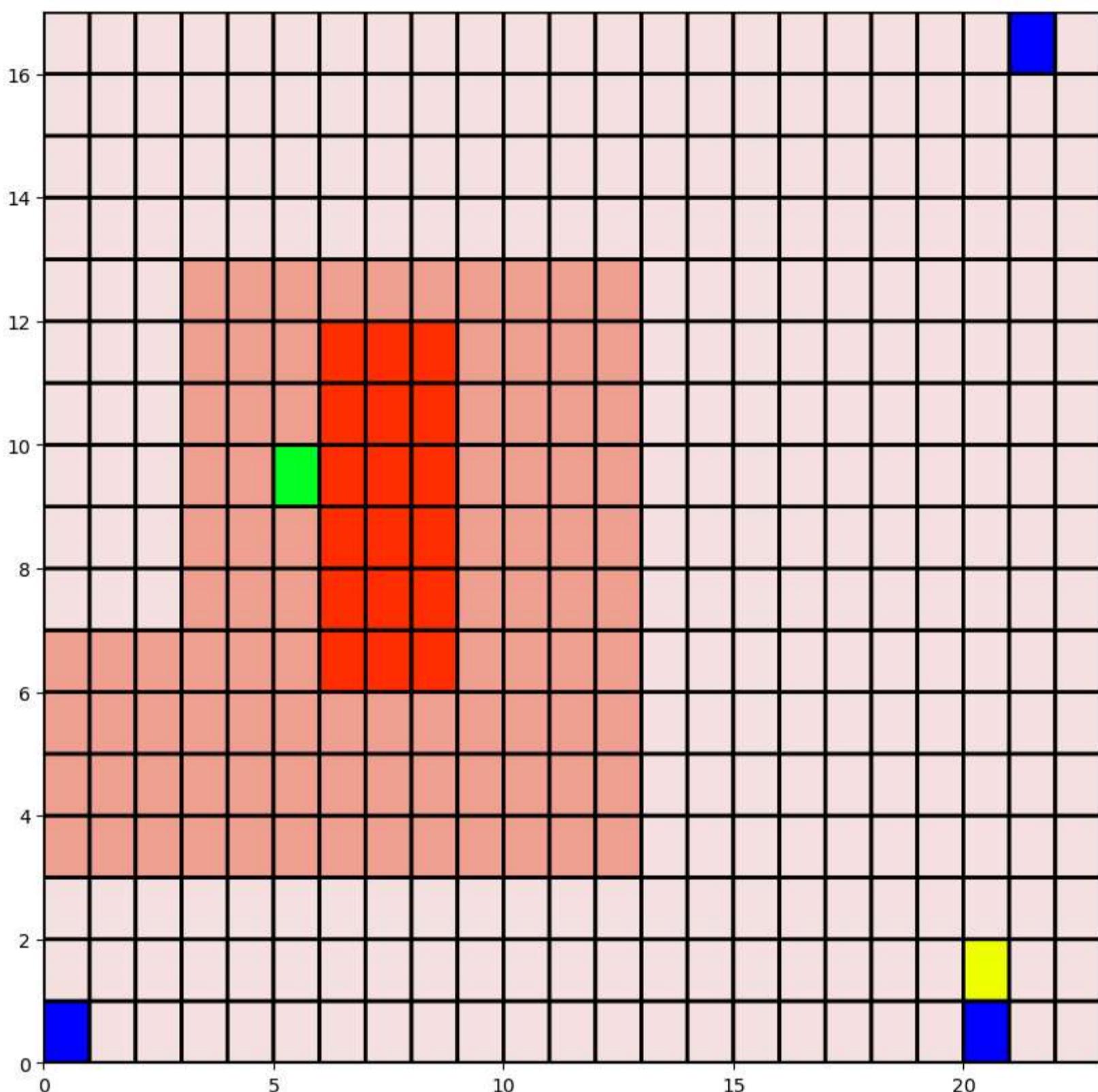
TODO: Repeat your experiments on the windy GridWorld

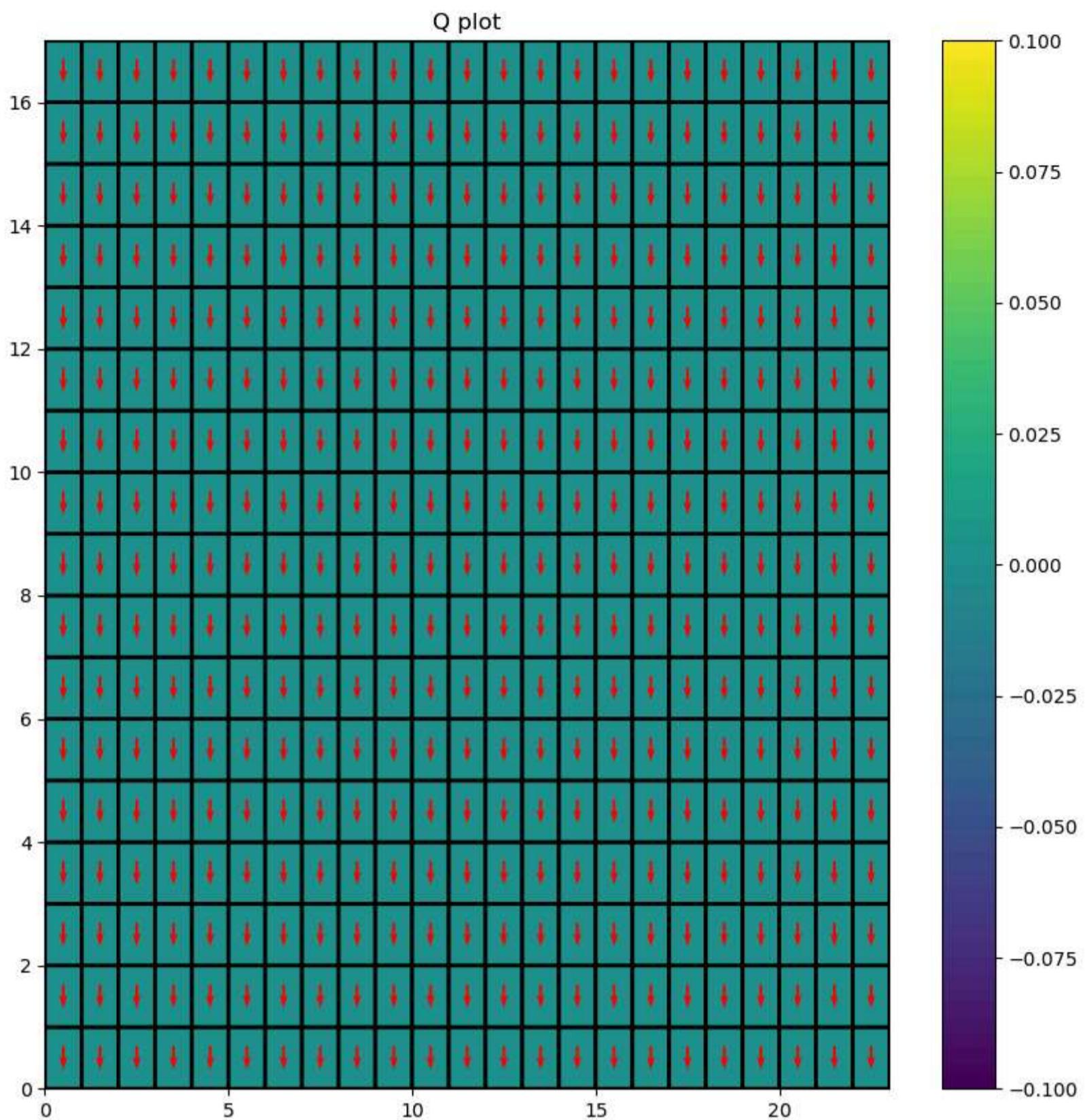
In [51]:

```
1 ### WRITE CODE HERE ### (set env to GridWorldWindyEnv)
2 world = 'grid_world2.txt'
3 goal_reward = 100
4 start_states = [(0,0), (0,20), (16,21)]
5 goal_states=[(9,5)]
6 max_steps=1000
7
8 from grid_world import GridWorldEnv, GridWorldWindyEnv
9 env = GridWorldWindyEnv(world, goal_reward=goal_reward, start_states=start_states, goal_states=goal_states,
10                         max_steps=max_steps, action_fail_prob=0.2)
11 plt.figure(figsize=(10, 10))
12 # Go UP
13 env.step(UP)
14 env.render(ax=plt, render_agent=True)
15
```

In [52]:

```
1 from grid_world import plot_Q
2
3 Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
4
5 plot_Q(Q)
6
7 Q.shape
```





Out[52]: (17, 23, 4)

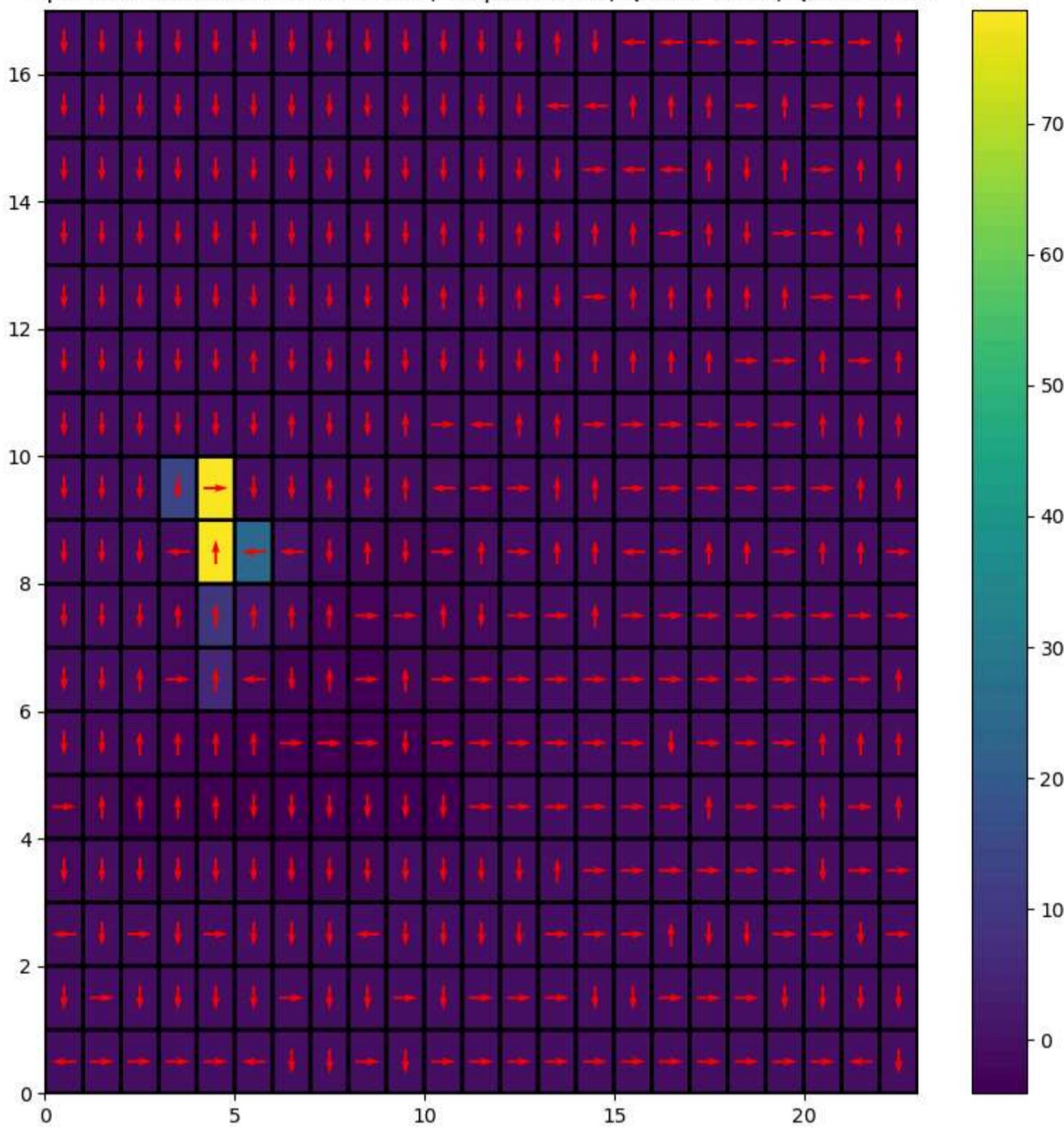
In [53]:

```
1 # initialize Q-value
2 Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
3
4 alpha0 = 0.4
5 gamma = 0.9
6 episodes = 5000
7 epsilon0 = 0.1
```

In [54]:

```
1 ### WRITE CODE HERE ### (repeat experiments - SARSA)
2 print_freq = 10
3
4 def sarsa(env, Q, gamma = 0.9, plot_heat = False, choose_action = choose_action_softmax):
5
6     episode_rewards = np.zeros(episodes)
7     steps_to_completion = np.zeros(episodes)
8     if plot_heat:
9         clear_output(wait=True)
10        plot_Q(Q)
11    epsilon = epsilon0
12    alpha = alpha0
13    for ep in tqdm(range(episodes)):
14        tot_reward, steps = 0, 0
15
16        # Reset environment
17        state = env.reset()
18        action = choose_action(Q, state)
19        done = False
20        while not done:
21            state_next, reward, done = env.step(action)
22            action_next = choose_action(Q, state_next)
23
24            # update equation
25            Q[state[0], state[1], action] += alpha*(reward + gamma*Q[state_next[0], state_next[1], action_next] -
26                                            Q[state[0], state[1], action])
27
28            tot_reward += reward
29            steps += 1
30
31            state, action = state_next, action_next
32
33        episode_rewards[ep] = tot_reward
34        steps_to_completion[ep] = steps
35
36        if (ep+1)%print_freq == 0 and plot_heat:
37            clear_output(wait=True)
38            plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f, Qmax: %.2f, Qmin: %.2f"%
39                                (ep+1, np.mean(episode_rewards[ep-print_freq+1:ep]),
40                                 np.mean(steps_to_completion[ep-print_freq+1:ep]),
41                                 Q.max(), Q.min()))
42
43    return Q, episode_rewards, steps_to_completion
44
45
46 Q, rewards, steps = sarsa(env, Q, gamma = gamma, plot_heat=True, choose_action=choose_action_softmax)
47
```

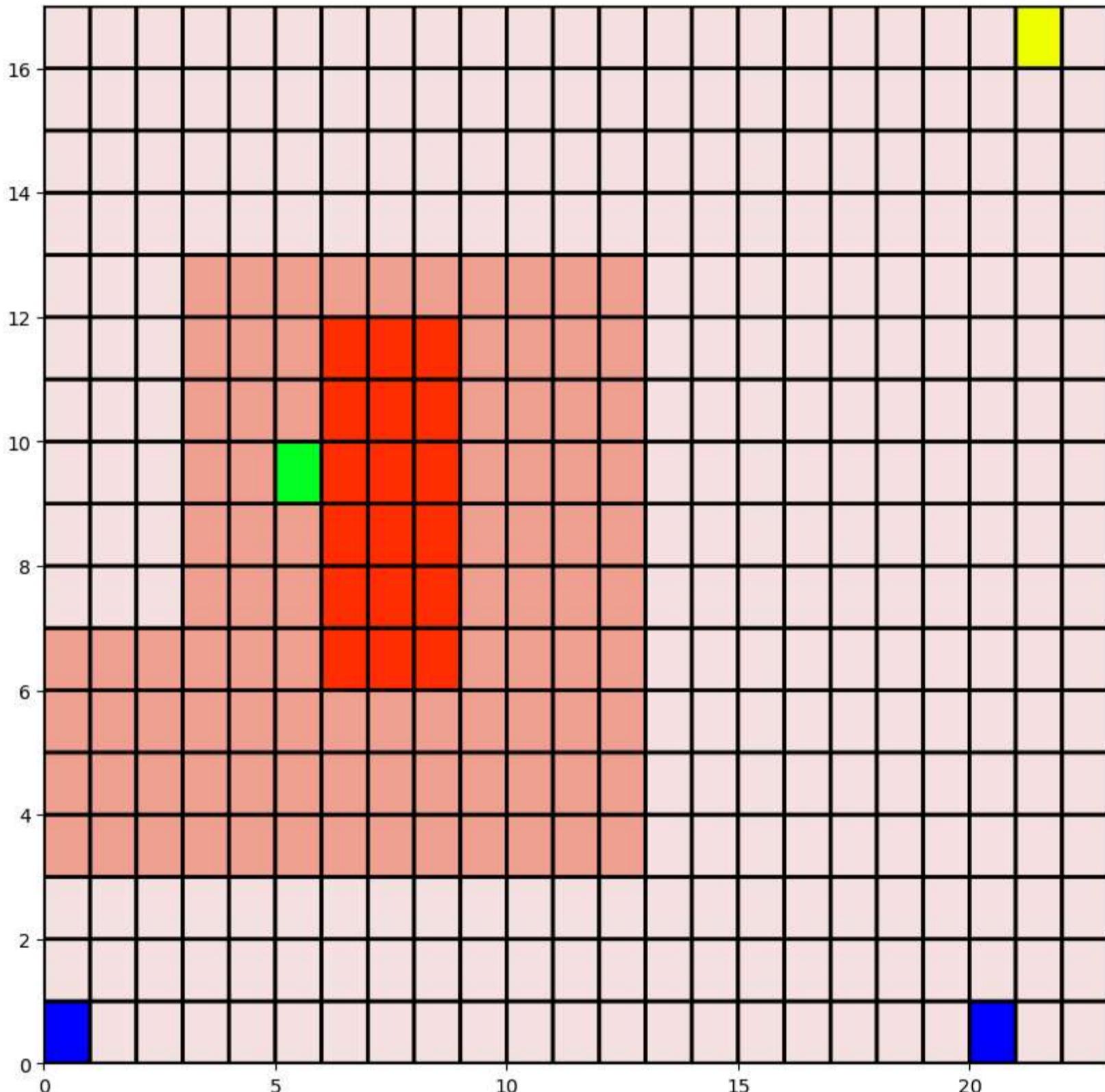
Episode 5000: Reward: -0.333333, Steps: 666.44, Qmax: 78.61, Qmin: -5.31



100% |  | 5000/5000 [02:02<00:00, 40.93it/s]

In [55]:

```
1 # Visualizing the policy Learnt by SARSA
2 from time import sleep
3
4 state = env.reset()
5 done = False
6 steps = 0
7 tot_reward = 0
8 while not done:
9     clear_output(wait=True)
10    state, reward, done = env.step(Q[state[0], state[1]].argmax())
11    plt.figure(figsize=(10, 10))
12    env.render(ax=plt, render_agent=True)
13    plt.show()
14    steps += 1
15    tot_reward += reward
16    sleep(0.2)
17 print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```



Steps: 659, Total Reward: 0

In [56]:

```
1 # Analyzing performance of the policy learnt by SARSA
2 Q_avgs_windy, reward_avgs_windy, steps_avgs_windy = [], [], []
3 num_expts = 3
4
5 for i in range(num_expts):
6     print("Experiment: %d" %(i+1))
7     Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
8     rg = np.random.RandomState(i)
9     Q, rewards, steps = sarsa(env, Q)
10    Q_avgs_windy.append(Q.copy())
11    reward_avgs_windy.append(rewards)
12    steps_avgs_windy.append(steps)
13
14 plt.xlabel('Episode')
15 plt.ylabel('Number of steps to Goal')
16 plt.title('SARSA: Steps to Goal')
17 plt.plot(np.arange(episodes),np.average(steps_avgs_windy, 0))
18 plt.show()
19 plt.xlabel('Episode')
20 plt.ylabel('Total Reward')
21 plt.title('SARSA: Total Reward')
22 plt.plot(np.arange(episodes),np.average(reward_avgs_windy, 0))
23 plt.show()
```

Experiment: 1

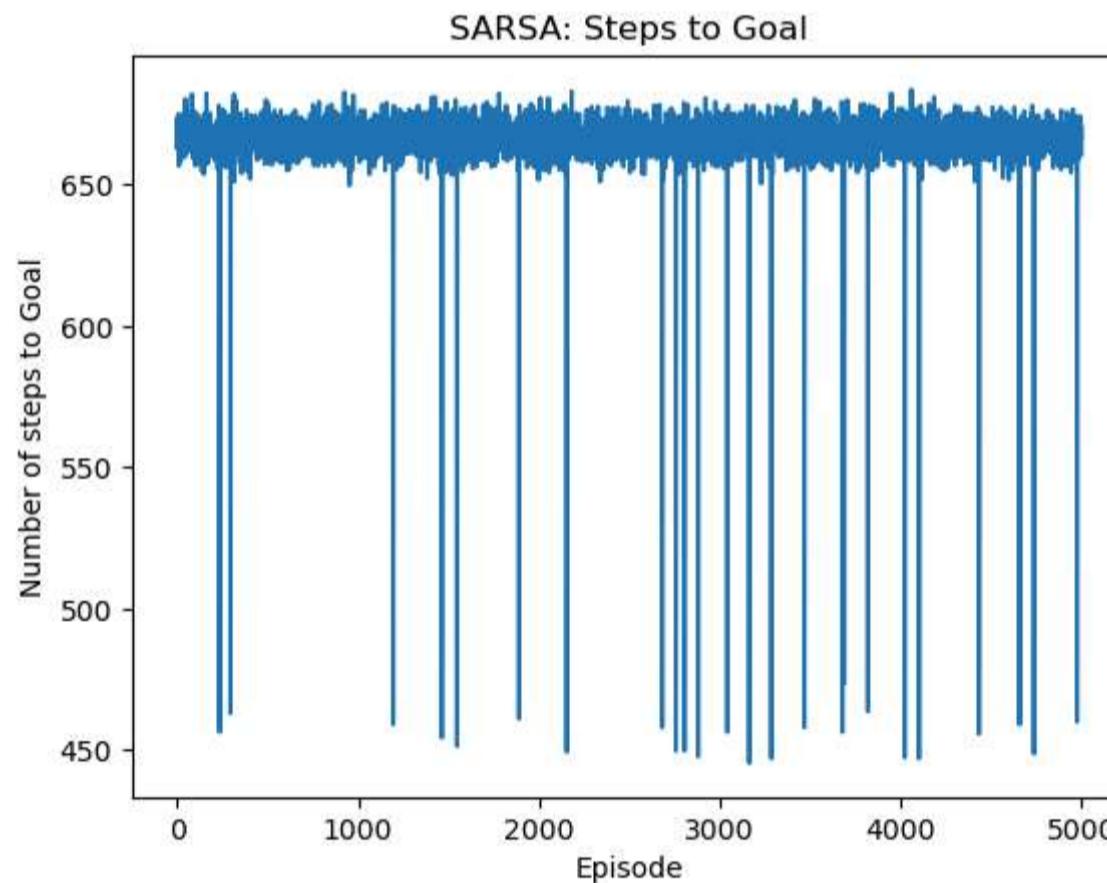
100%|██████████| 5000/5000 [01:06<00:00, 75.49it/s]

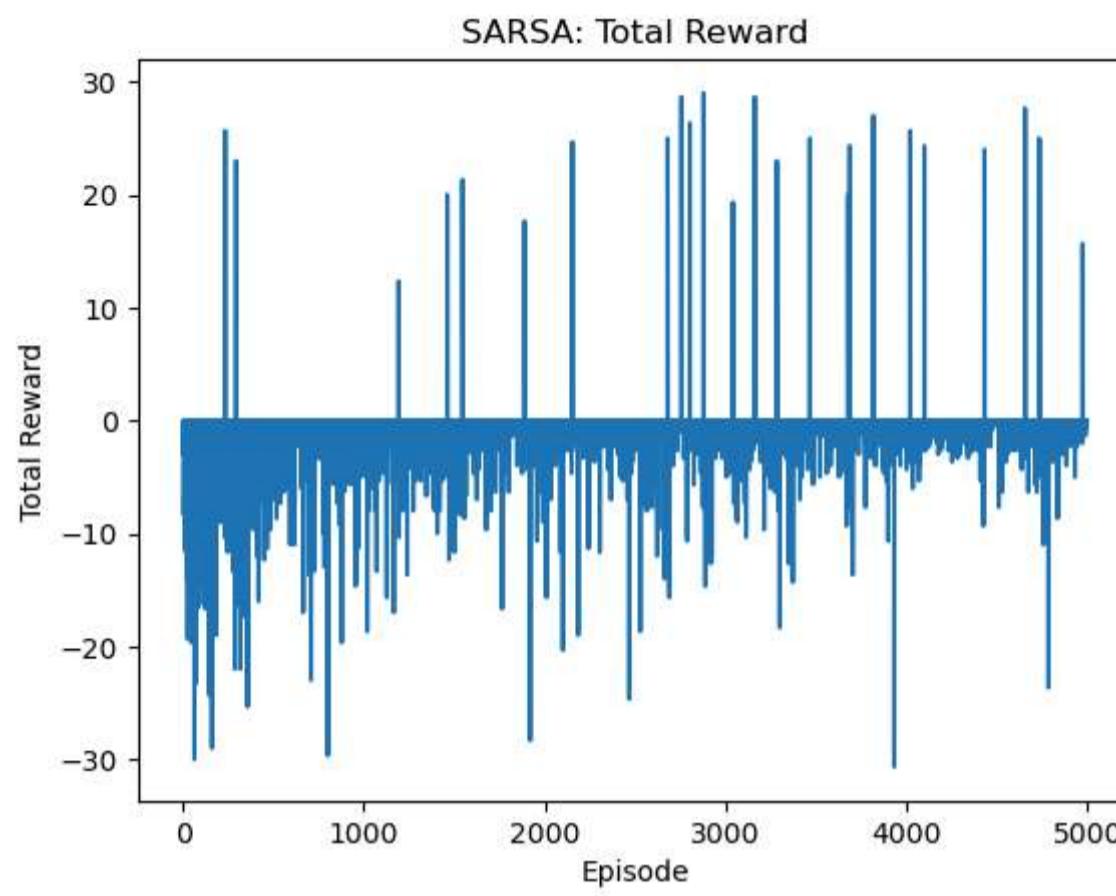
Experiment: 2

100%|██████████| 5000/5000 [01:06<00:00, 75.29it/s]

Experiment: 3

100%|██████████| 5000/5000 [01:07<00:00, 74.54it/s]

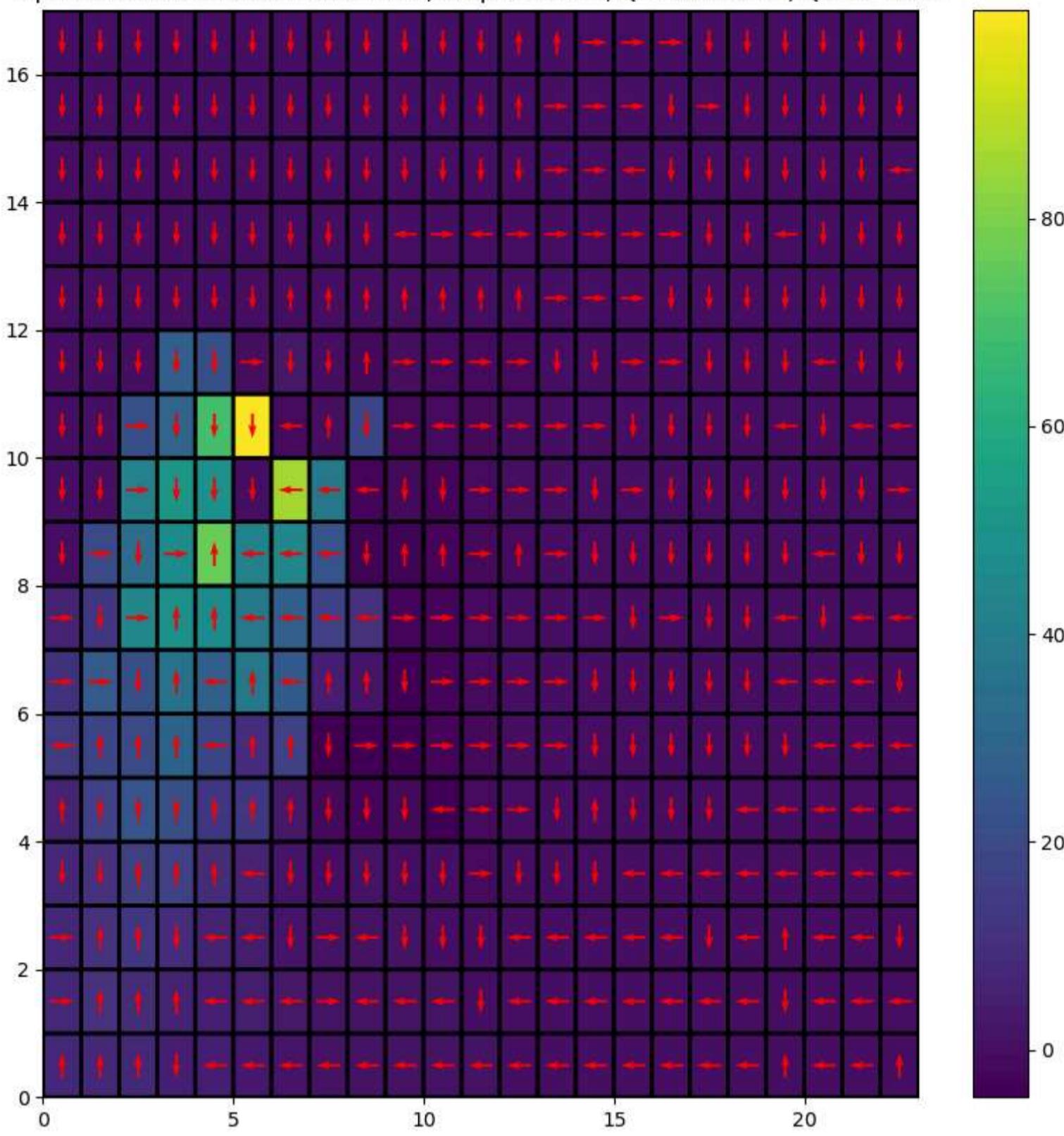




In [66]:

```
1 ### WRITE CODE HERE ### (repeat experiments - Q-Learning)
2 print_freq = 10
3
4 def QLearning(env, Q, gamma = 0.9, plot_heat = False, choose_action = choose_action_softmax):
5
6     episode_rewards = np.zeros(episodes)
7     steps_to_completion = np.zeros(episodes)
8     if plot_heat:
9         clear_output(wait=True)
10        plot_Q(Q)
11
12    epsilon = epsilon0
13    alpha = alpha0
14    for ep in tqdm(range(episodes)):
15        tot_reward, steps = 0, 0
16
17        state = env.reset()
18        done = False
19        while not done:
20            action = choose_action(Q, state) # -> softmax
21            #action = choose_action(Q, state, epsilon) # -> Epsilon greedy
22            state_next, reward, done = env.step(action)
23
24            Q[state[0], state[1], action] += alpha*(reward + gamma*np.max(Q[state_next[0], state_next[1], :]) -
25                                              Q[state[0], state[1], action])
26            tot_reward += reward
27            steps += 1
28
29            state = state_next
30
31        episode_rewards[ep] = tot_reward
32        steps_to_completion[ep] = steps
33
34        if (ep+1)%print_freq == 0 and plot_heat:
35            clear_output(wait=True)
36            plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f, Qmax: %.2f, Qmin: %.2f"%
37                                         (ep+1, np.mean(episode_rewards[ep-print_freq+1:ep]),
38                                          np.mean(steps_to_completion[ep-print_freq+1:ep]),
39                                          Q.max(), Q.min()))
39
40    return Q, episode_rewards, steps_to_completion
41
42
43 Q, rewards, steps = QLearning(env, Q, gamma = gamma, plot_heat=True, choose_action=choose_action_softmax)
```

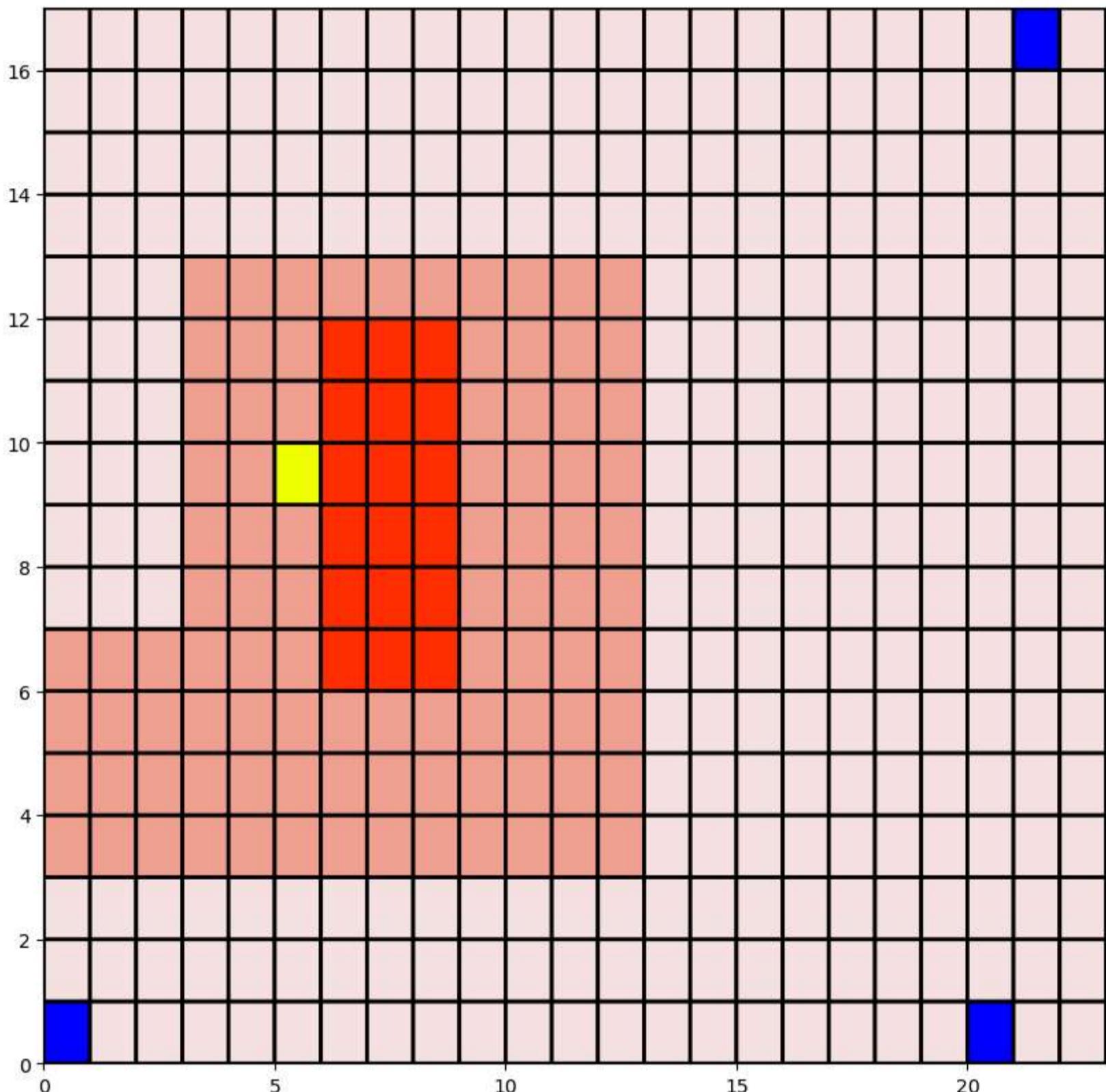
Episode 5000: Reward: 14.666667, Steps: 528.33, Qmax: 100.00, Qmin: -5.83



100% | | 5000/5000 [01:56<00:00, 42.94it/s]

In [69]:

```
1 # Visualizing the policy Learnt by QLearning
2 from time import sleep
3
4 state = env.reset()
5 done = False
6 steps = 0
7 tot_reward = 0
8 while not done:
9     clear_output(wait=True)
10    state, reward, done = env.step(Q[state[0], state[1]].argmax())
11    plt.figure(figsize=(10, 10))
12    env.render(ax=plt, render_agent=True)
13    plt.show()
14    steps += 1
15    tot_reward += reward
16    sleep(0.2)
17 print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```



Steps: 215, Total Reward: 88

In [59]:

```
1 # Analyzing performance of the policy learnt by QLearning
2 Q_avgs_Qlearning_windy, reward_avgs_Qlearning_windy, steps_avgs_Qlearning_windy = [], [], []
3 num_expts = 3
4
5 for i in range(num_expts):
6     print("Experiment: %d" %(i+1))
7     Q = np.zeros((env.grid.shape[0], env.grid.shape[1], len(env.action_space)))
8     rg = np.random.RandomState(i)
9     Q, rewards, steps = sarsa(env, Q)
10    Q_avgs_Qlearning_windy.append(Q.copy())
11    reward_avgs_Qlearning_windy.append(rewards)
12    steps_avgs_Qlearning_windy.append(steps)
13
14 plt.xlabel('Episode')
15 plt.ylabel('Number of steps to Goal')
16 plt.title('Q-Learning: Steps to Goal')
17 plt.plot(np.arange(episodes), np.average(steps_avgs_Qlearning_windy, 0), color='red')
18 plt.show()
19
20 plt.xlabel('Episode')
21 plt.ylabel('Total Reward')
22 plt.title('Q-Learning: Total Reward')
23 plt.plot(np.arange(episodes), np.average(reward_avgs_Qlearning_windy, 0), color='red')
24 plt.show()
```

Experiment: 1

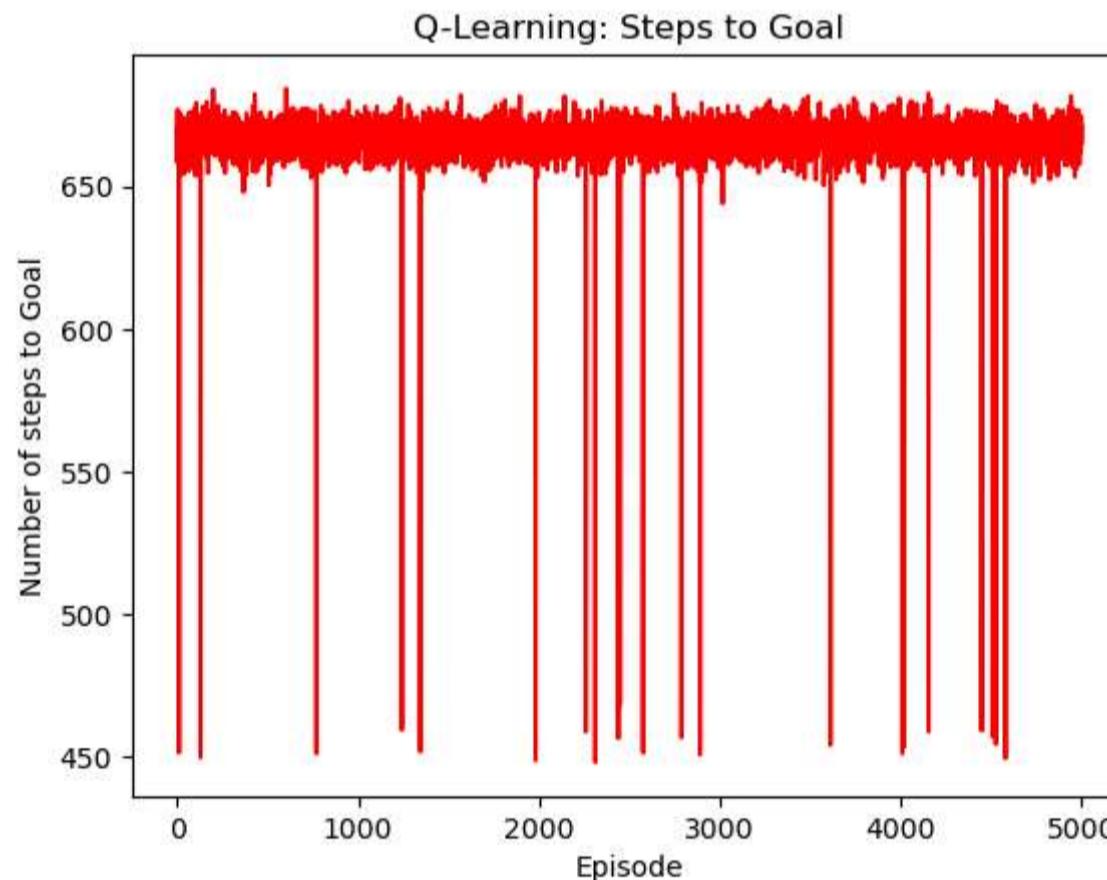
100%|██████████| 5000/5000 [01:05<00:00, 75.99it/s]

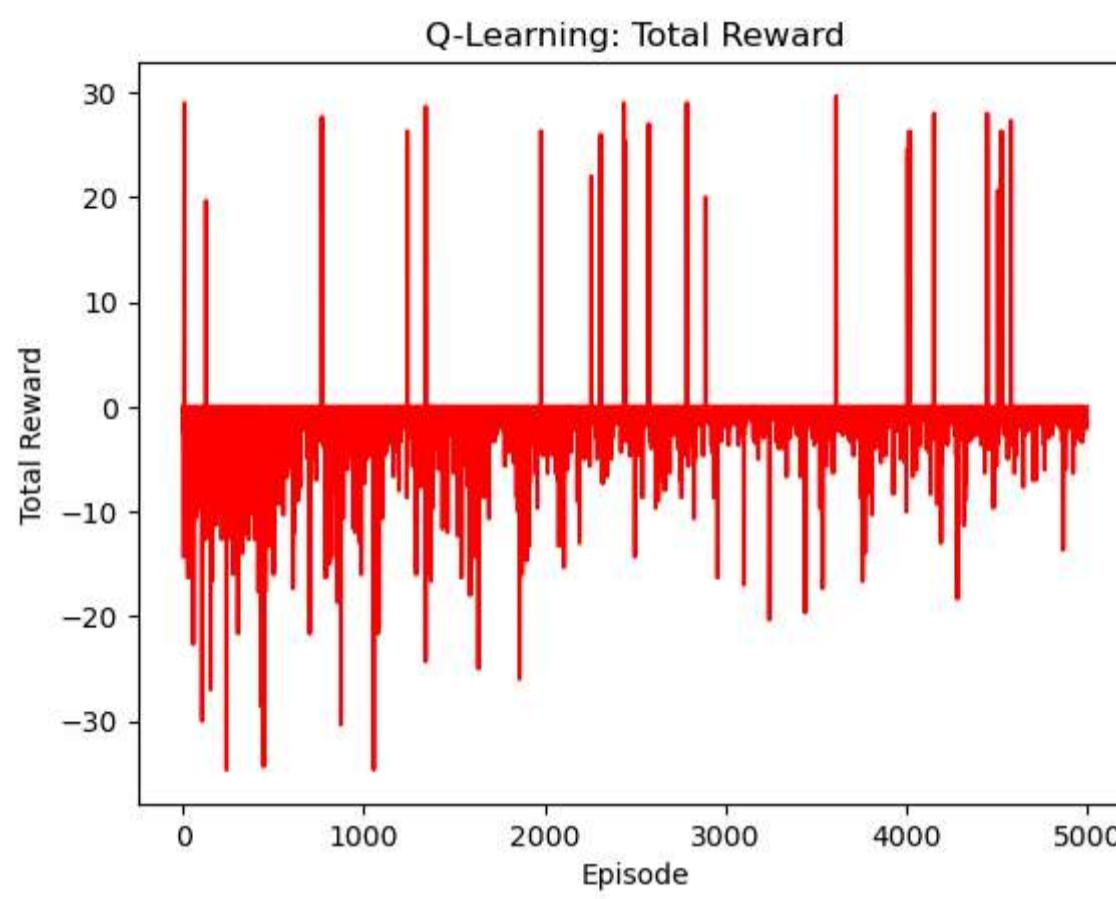
Experiment: 2

100%|██████████| 5000/5000 [01:08<00:00, 72.74it/s]

Experiment: 3

100%|██████████| 5000/5000 [01:06<00:00, 75.26it/s]





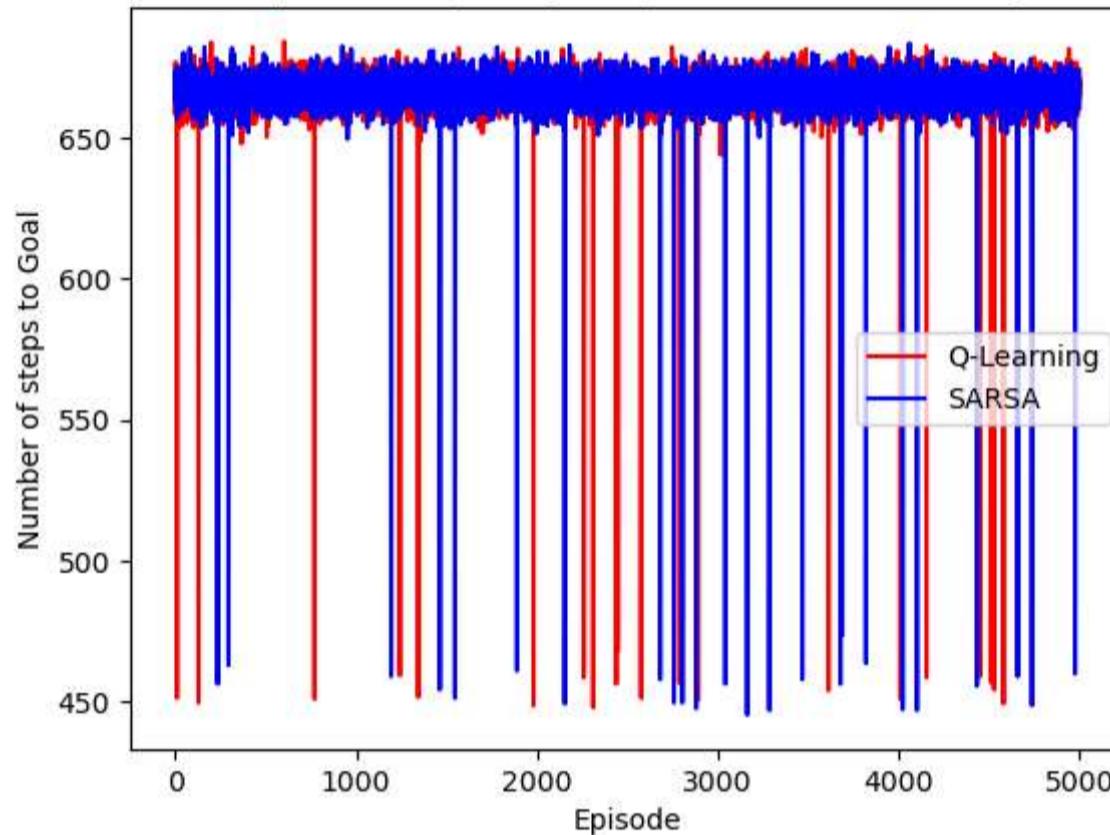
In [60]:

```

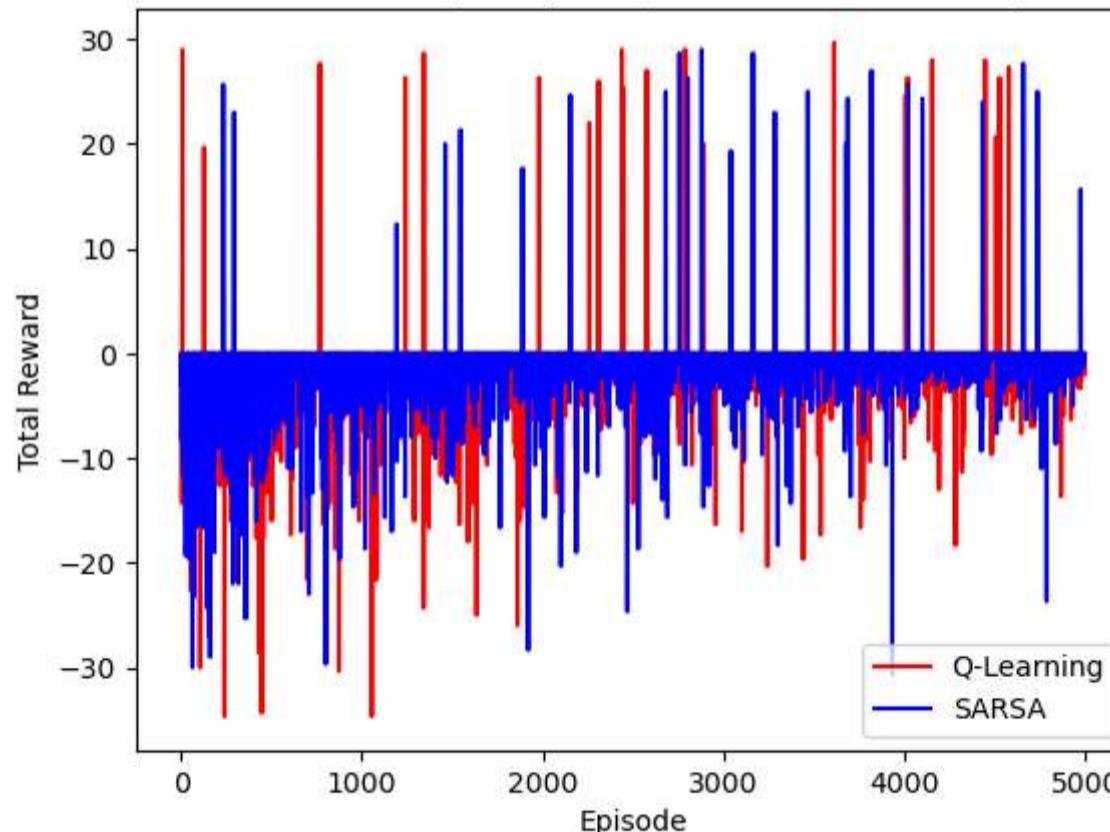
1 ### WRITE CODE HERE ### (compare Q-Learning and SARSA performance)
2
3 plt.xlabel('Episode')
4 plt.ylabel('Number of steps to Goal')
5 plt.title('Steps to Goal(Windy Env): SARSA vs Q-Learning')
6 plt.plot(np.arange(episodes), np.average(steps_avgs_Qlearning_windy, 0), color='red', label='Q-Learning')
7 plt.plot(np.arange(episodes), np.average(steps_avgs_windy, 0), color='blue', label='SARSA')
8 plt.legend()
9 plt.show()
10
11
12 plt.xlabel('Episode')
13 plt.ylabel('Total Reward')
14 plt.title('Total Reward(Windy Env): SARSA vs Q-Learning')
15 plt.plot(np.arange(episodes), np.average(reward_avgs_Qlearning_windy, 0), color='red', label='Q-Learning')
16 plt.plot(np.arange(episodes), np.average(reward_avgs_windy, 0), color='blue', label='SARSA')
17 plt.legend()
18 plt.show()
19

```

Steps to Goal(Windy Env): SARSA vs Q-Learning



Total Reward(Windy Env): SARSA vs Q-Learning



For the Windy Environment:

1. Sarsa Algorithm doesn't learn a policy to reach the goal state as it tends to optimize the policy according to the behaviour policy(softmax/ epsilon greedy).
2. Q learning algorithm learns a policy to reach the goal state as it always updates Q-values using the action that maximizes future rewards, regardless of the behavior policy. Due to this Off policy nature the QLearning algo is able to find a path to the goal state with considerable rewards value in presence of the Windy nature.

The Average Reward of Q learning is Higher than the SARSA algorithm. The Q learning algorithm has learnt a policy to reach the goal state from all the three start positions with considerable reward, whereas the agent following the SARSA policy gets stuck at the right most cells due the windy nature of the environment and poor policy learnt.

1