

CS6700: Reinforcement Learning - Tutorial 4 (Q-Learning and SARSA)

Your tasks are as follows:

1. Complete code for ϵ -greedy and softmax action selection policy
2. Complete update equation for SARSA - train and visualize an agent
3. Analyze performance of SARSA - Plot total reward & steps taken per episode (averaged across 5 runs)
4. Complete update equation for Q-Learning - train and visualize an agent
5. Analyze performance of Q-Learning - Plot total reward & steps taken per episode (averaged across 5 runs)

```
import numpy as np
import matplotlib.pyplot as plt
from tqdm import tqdm
from IPython.display import clear_output
%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.

```
DOWN = 0
UP = 1
LEFT = 2
RIGHT = 3
actions = [DOWN, UP, LEFT, RIGHT]
```

Let us construct a grid in a text file.

```
!cat 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
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 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 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.

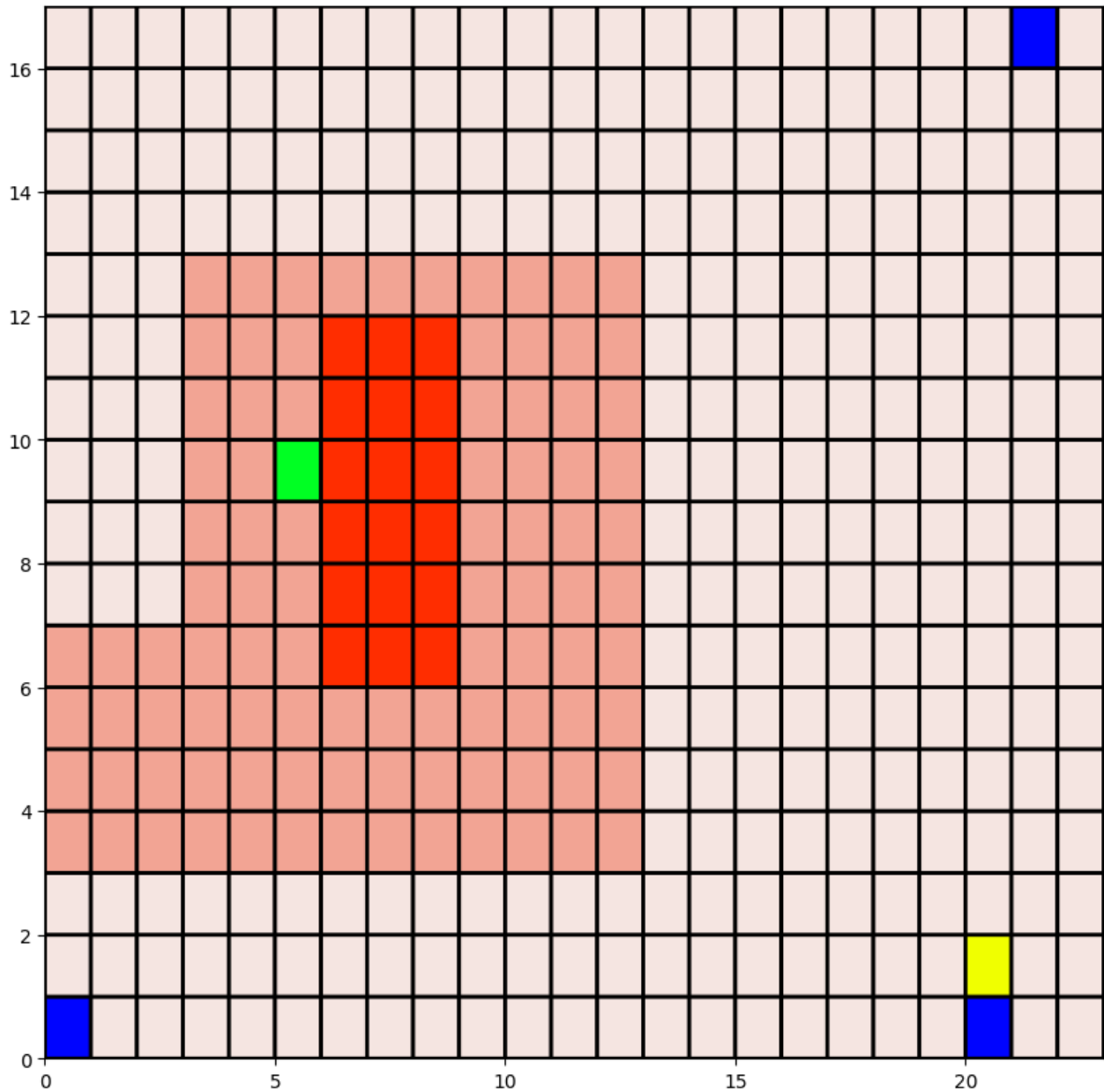
```

world = 'grid_world2.txt'
goal_reward = 100
start_states = [(0,0), (0,20), (16,21)]
goal_states=[(9,5)]
max_steps=10000

from grid_world import GridWorldEnv, GridWorldWindyEnv

env = GridWorldEnv(world, goal_reward=goal_reward,
start_states=start_states, goal_states=goal_states,
max_steps=max_steps, action_fail_prob=0.2)
plt.figure(figsize=(10, 10))
# Go UP
env.step(UP)
env.render(ax=plt, render_agent=True)

```



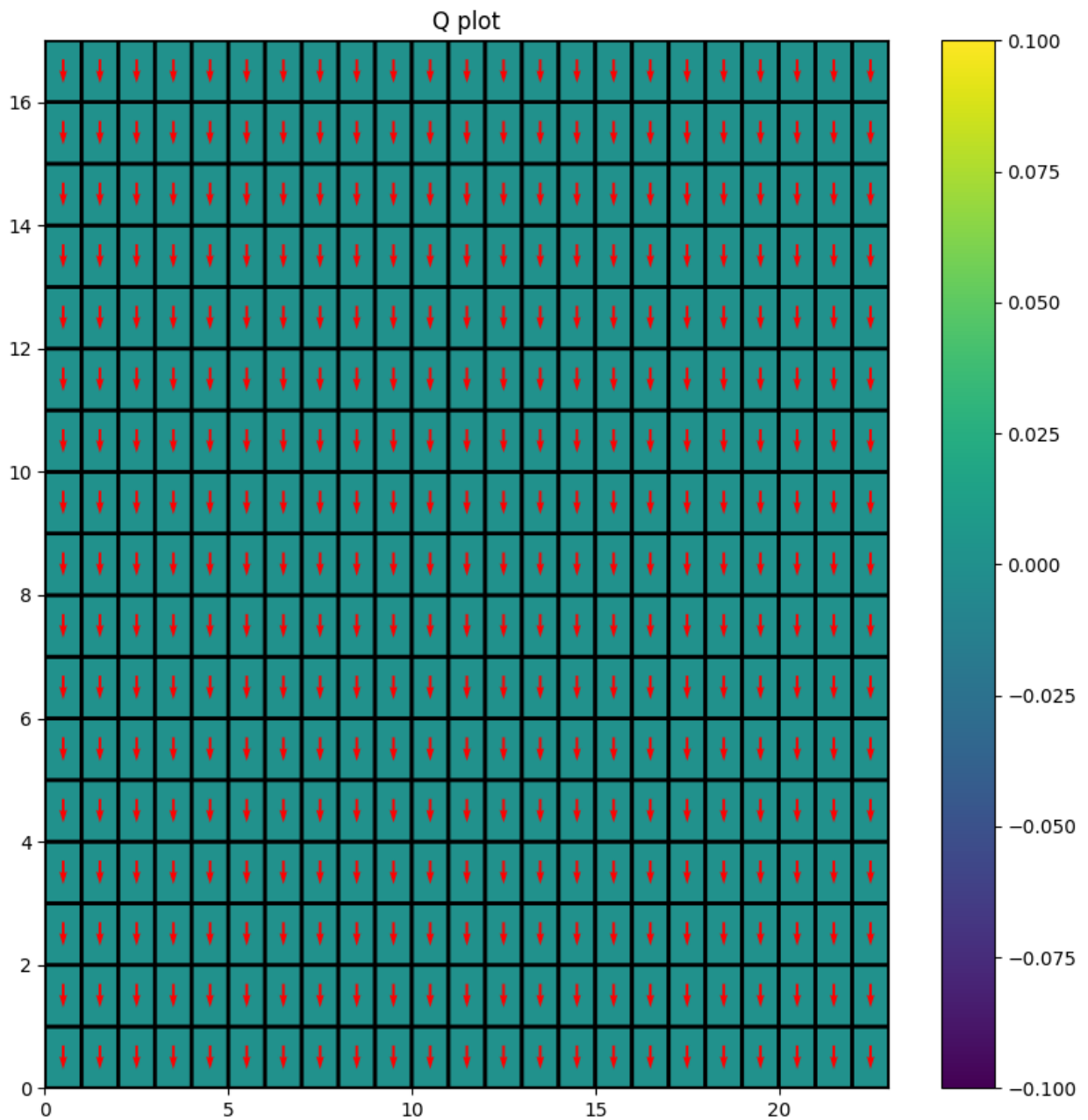
Legend

- *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.

```
from grid_world import plot_Q
Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
len(env.action_space)))
plot_Q(Q)
Q.shape
```



```
(17, 23, 4)
```

Exploration strategies

1. Epsilon-greedy
2. Softmax

```
from scipy.special import softmax

seed = 42
rg = np.random.RandomState(seed)

# Epsilon greedy
def choose_action_epsilon(Q, state, epsilon, rg=rg):
    if rg.rand() < epsilon or not Q[state[0], state[1]].any(): # TODO:
        eps greedy condition
        return np.random.randint(0, len(actions)) # TODO: return
        random action
    else:
        Q_actions = Q[state[0], state[1], :]
        return np.argmax((Q_actions)) # TODO: return best action

# Softmax
def choose_action_softmax(Q, state, _, rg=rg):
    Q_actions = np.array(Q[state[0], state[1], :])
    probs = softmax(Q_actions)
    action = np.random.choice([i for i in range(len(actions))],
                              p=probs)
    return action # TODO: return random action with selection
    probability
```

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

```
# initialize Q-value
Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
              len(env.action_space)))

alpha0 = 0.4
gamma = 0.9
```

```
episodes = 10000
epsilon0 = 0.1
```

Let's implement SARSA

```
print_freq = 100

def sarsa(env, Q, gamma = 0.9, plot_heat = False, choose_action =
choose_action_softmax):

    episode_rewards = np.zeros(episodes)
    steps_to_completion = np.zeros(episodes)
    if plot_heat:
        clear_output(wait=True)
        plot_Q(Q)
    epsilon = epsilon0
    alpha = alpha0
    for ep in tqdm(range(episodes)):
        tot_reward, steps = 0, 0

        # Reset environment
        state = env.reset()
        action = choose_action(Q, state, epsilon)
        done = False
        while not done:
            state_next, reward, done = env.step(action)
            action_next = choose_action(Q, state_next, epsilon)

            # TODO: update equation
            Q[state[0], state[1], action] = Q[state[0], state[1],
action] + alpha * (reward + gamma * Q[state_next[0], state_next[1],
action_next] - Q[state[0], state[1], action])

            tot_reward += reward
            steps += 1

            state, action = state_next, action_next

        episode_rewards[ep] = tot_reward
        steps_to_completion[ep] = steps

        if (ep+1)%print_freq == 0 and plot_heat:
            clear_output(wait=True)
            plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f,
Qmax: %.2f, Qmin: %.2f"%(ep+1, np.mean(episode_rewards[ep-
print_freq+1:ep]),
np.mean(steps_to_completion[ep-print_freq+1:ep]),
```

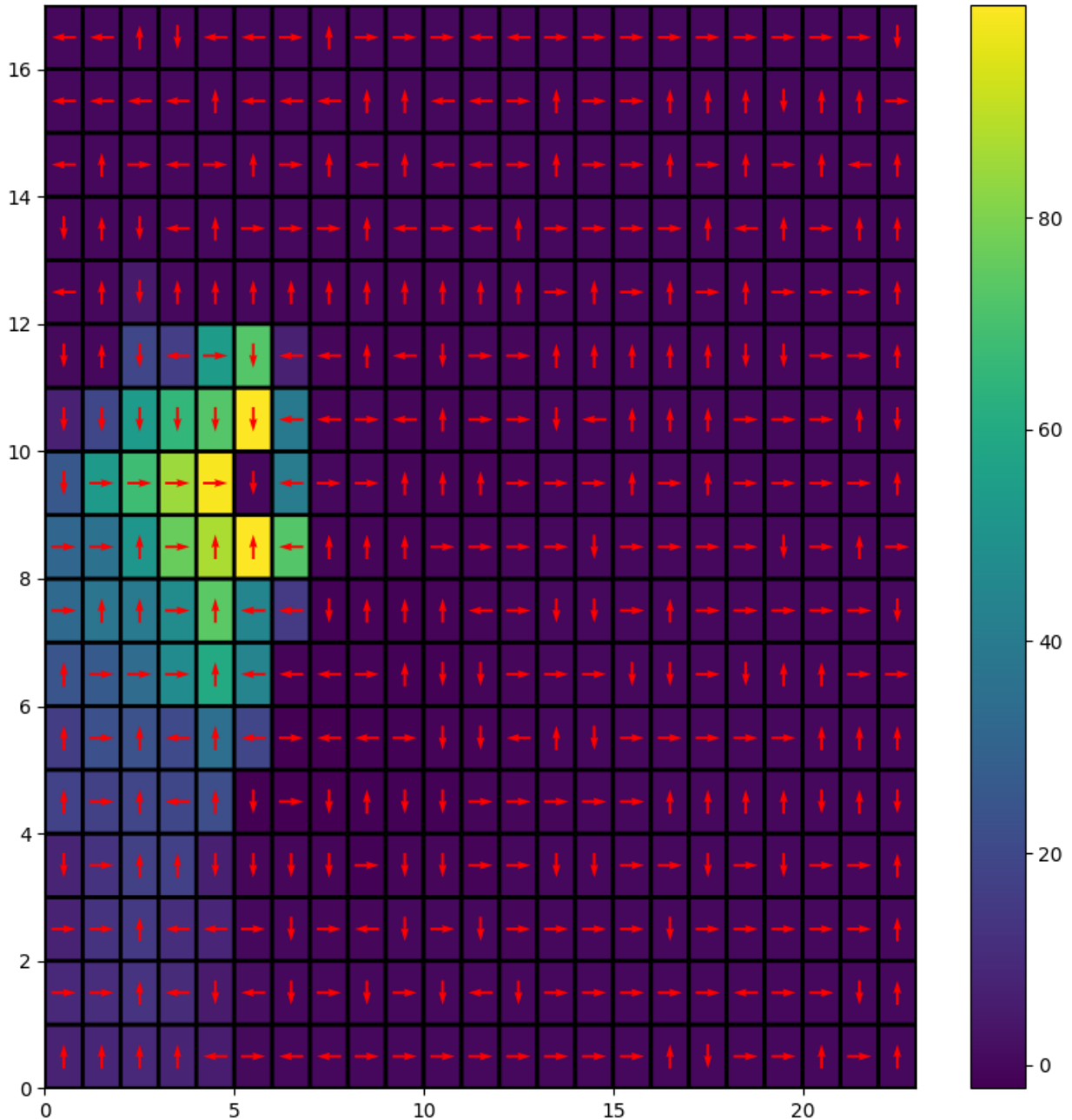
```
Q.max(), Q.min()))
```

```
return Q, episode_rewards, steps_to_completion
```

SARSA with ϵ Greedy Policy

```
Q, rewards, steps = sarsa(env, Q, gamma = gamma, plot_heat=True,  
choose_action= choose_action_epsilon)
```

Episode 10000: Reward: 29.848485, Steps: 6573.19, Qmax: 100.00, Qmin: -2.91



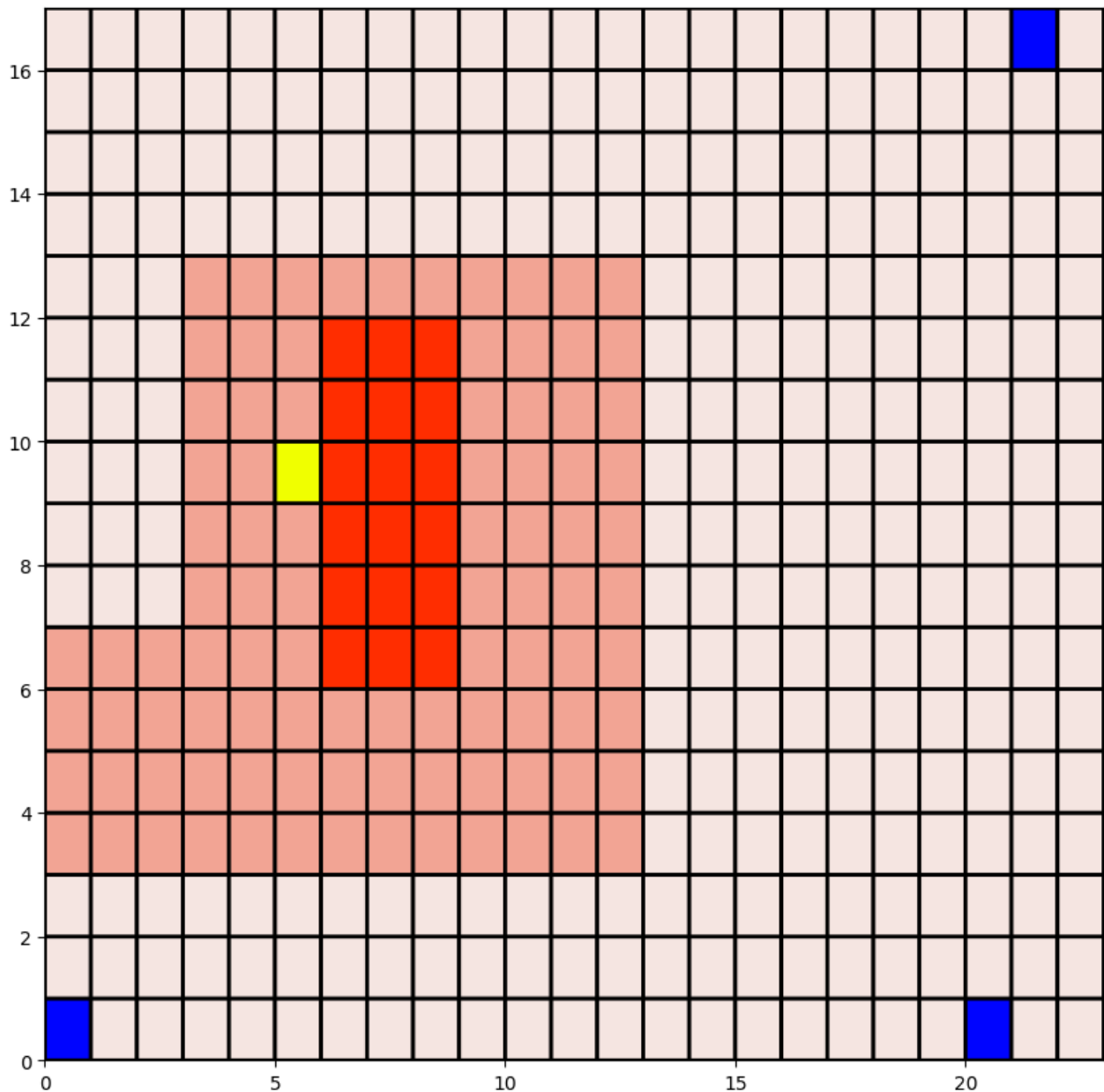
100%|██████████| 10000/10000 [12:28<00:00, 13.35it/s]

Visualizing the policy

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

```
from time import sleep

state = env.reset()
done = False
steps = 0
tot_reward = 0
while not done:
    clear_output(wait=True)
    state, reward, done = env.step(Q[state[0], state[1]].argmax())
    plt.figure(figsize=(10, 10))
    env.render(ax=plt, render_agent=True)
    plt.show()
    steps += 1
    tot_reward += reward
    sleep(0.2)
print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```

Steps: 19, Total Reward: 89

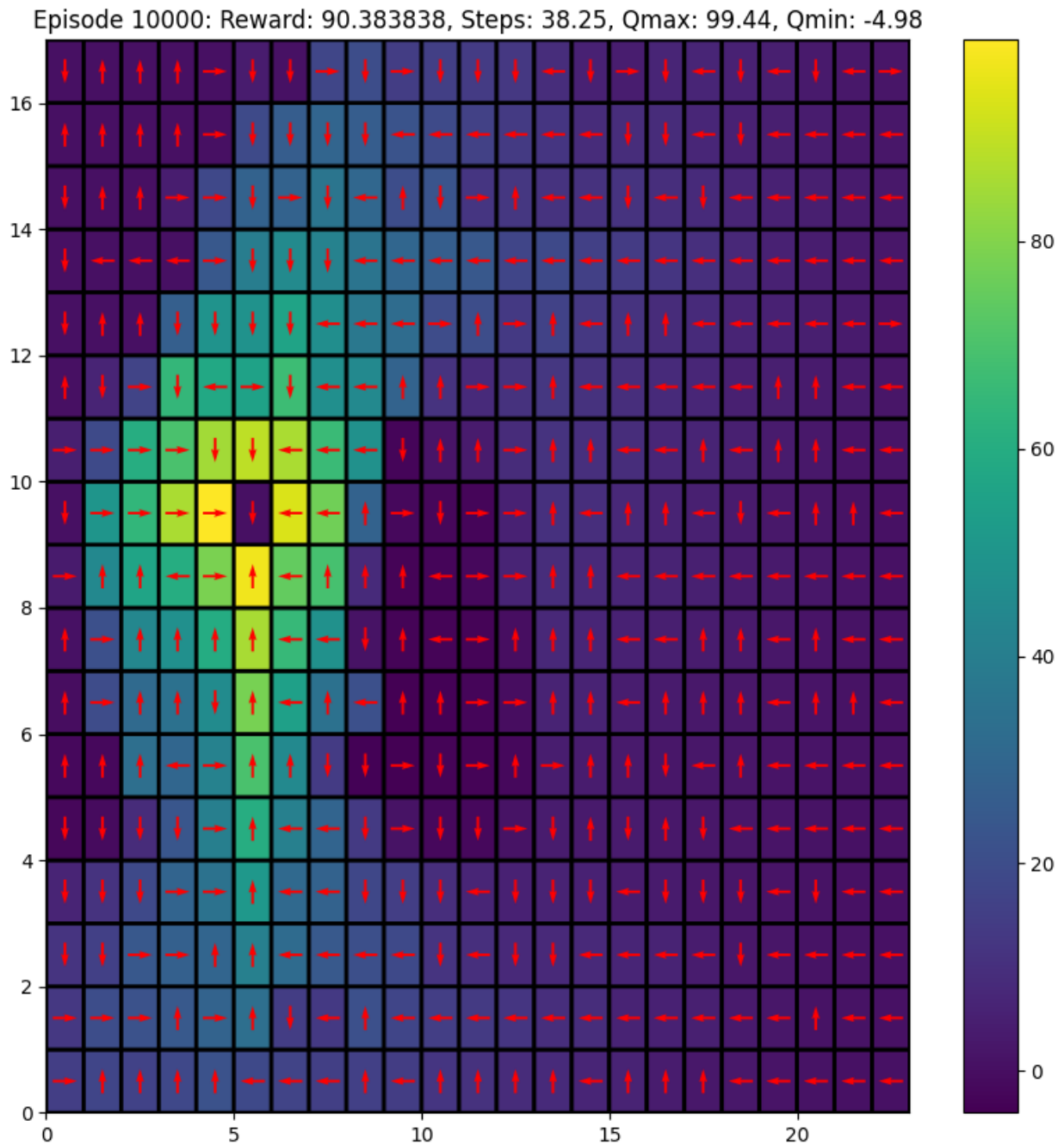
SARSA with Softmax Policy

Initializing Q values and Hyperparameters again

```
Q = np.zeros((env.grid.shape[0], env.grid.shape[1],  
len(env.action_space)))
```

```
alpha0 = 0.4  
gamma = 0.9  
episodes = 10000  
epsilon0 = 0.1
```

```
Q, rewards, steps = sarsa(env, Q, gamma = gamma, plot_heat=True,  
choose_action= choose_action_softmax)
```



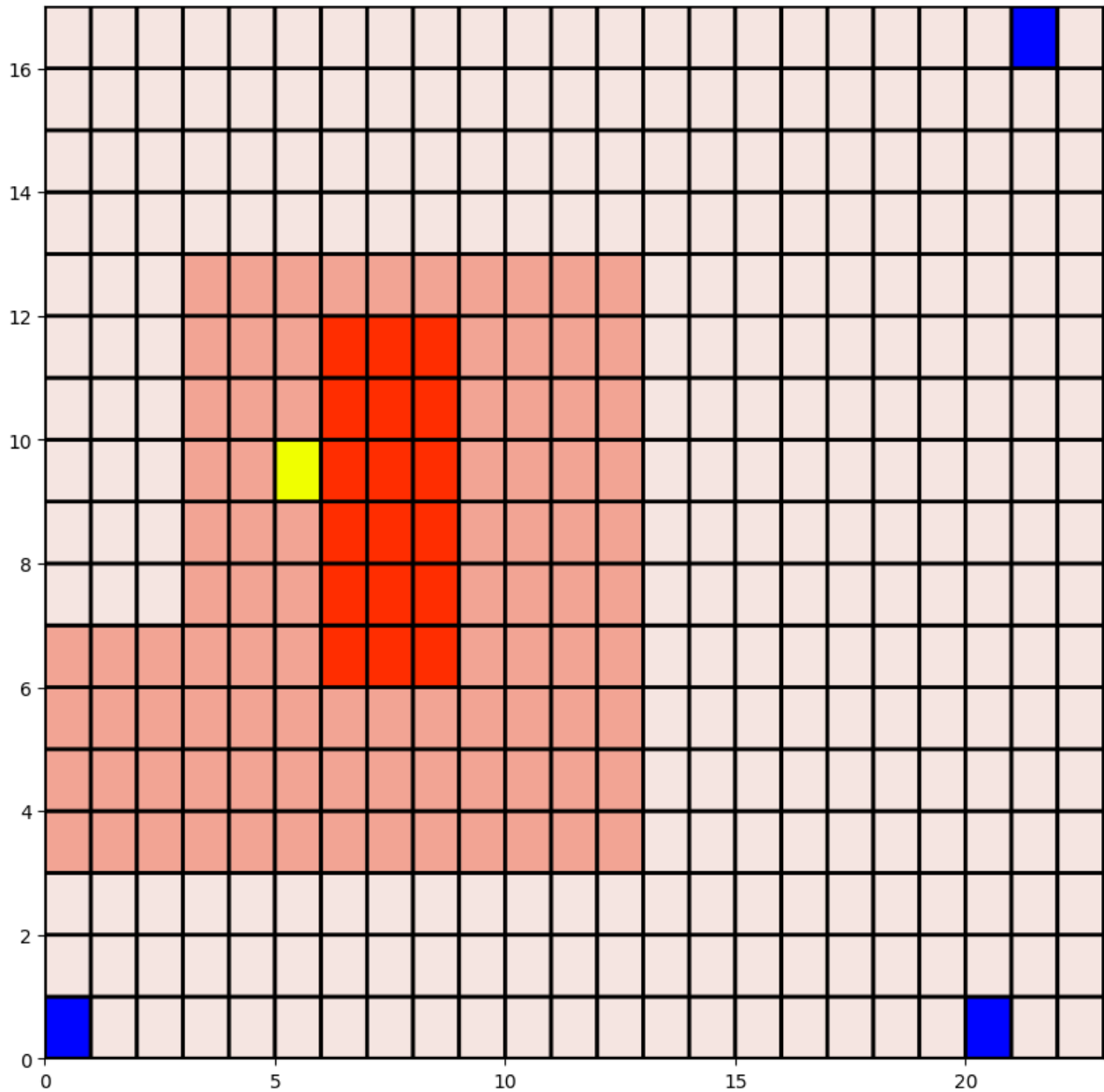
100%|██████████| 10000/10000 [00:39<00:00, 256.39it/s]

Visualizing the policy

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

```
from time import sleep

state = env.reset()
done = False
steps = 0
tot_reward = 0
while not done:
    clear_output(wait=True)
    state, reward, done = env.step(Q[state[0], state[1]].argmax())
    plt.figure(figsize=(10, 10))
    env.render(ax=plt, render_agent=True)
    plt.show()
    steps += 1
    tot_reward += reward
    sleep(0.2)
print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```



Steps: 30, Total Reward: 94

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.

```

num_expts = 5
reward_avgs, steps_avgs = [], []

for i in range(num_expts):
    print("Experiment: %d"%(i+1))
    Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
len(env.action_space)))
    rg = np.random.RandomState(i)

    # TODO: run sarsa, store metrics
    Q, rewards, steps = sarsa(env, Q, gamma = gamma, plot_heat=False,
choose_action = choose_action_softmax)
    reward_avgs.append(rewards)
    steps_avgs.append(steps)

reward_avgs = np.array(reward_avgs)
steps_avgs = np.array(steps_avgs)
reward_avgs = np.mean(reward_avgs, axis=0)
steps_avgs = np.mean(steps_avgs, axis=0)

Experiment: 1
100%|██████████| 10000/10000 [00:23<00:00, 417.81it/s]
Experiment: 2
100%|██████████| 10000/10000 [00:20<00:00, 481.24it/s]
Experiment: 3
100%|██████████| 10000/10000 [00:19<00:00, 518.26it/s]
Experiment: 4
100%|██████████| 10000/10000 [00:26<00:00, 379.08it/s]
Experiment: 5
100%|██████████| 10000/10000 [00:26<00:00, 375.16it/s]

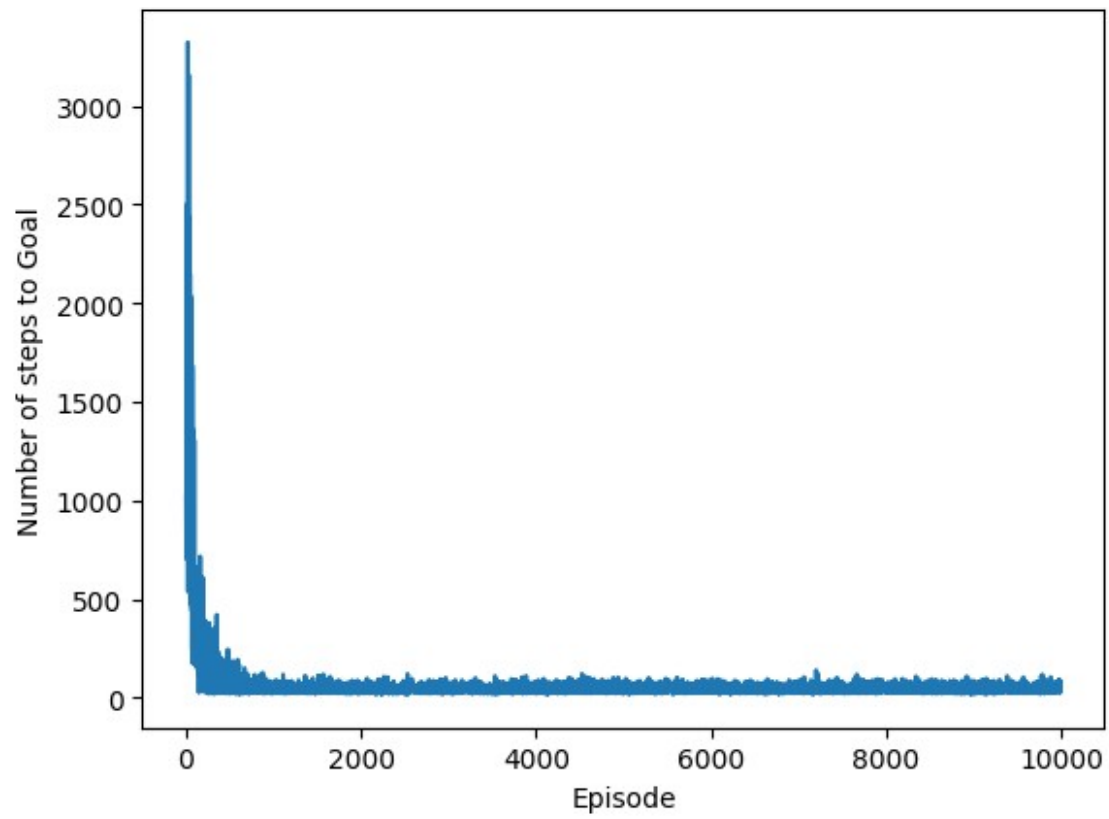
# TODO: visualize individual metrics vs episode count (averaged across
multiple run(s))

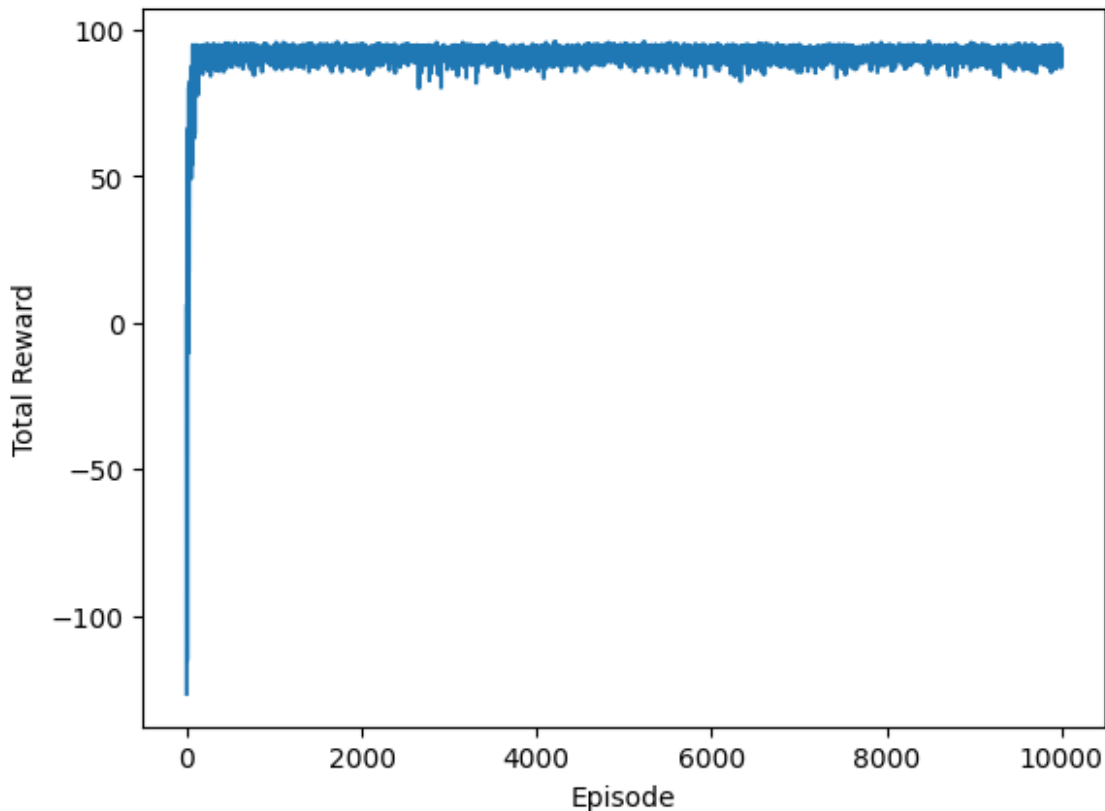
plt.figure()
plt.plot(steps_avgs)
plt.xlabel('Episode')
plt.ylabel('Number of steps to Goal')
plt.show()

plt.figure()
plt.plot(reward_avgs)
plt.xlabel('Episode')

```

```
plt.ylabel('Total Reward')  
plt.show()
```





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.

```
# initialize Q-value
Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
len(env.action_space)))

alpha0 = 0.4
gamma = 0.9
episodes = 10000
epsilon0 = 0.1

print_freq = 100

def qlearning(env, Q, gamma = 0.9, plot_heat = False, choose_action =
choose_action_softmax):

    episode_rewards = np.zeros(episodes)
    steps_to_completion = np.zeros(episodes)
```

```

if plot_heat:
    clear_output(wait=True)
    plot_Q(Q)
epsilon = epsilon0
alpha = alpha0
for ep in tqdm(range(epochs)):
    tot_reward, steps = 0, 0

    # Reset environment
    state = env.reset()
    action = choose_action(Q, state, epsilon)
    done = False
    while not done:
        state_next, reward, done = env.step(action)
        action_next = choose_action(Q, state_next, epsilon)

        # TODO: update equation
        Q[state[0], state[1], action] = Q[state[0], state[1],
action] + alpha * (reward + gamma * max(Q[state_next[0],
state_next[1], :]) - Q[state[0], state[1], action])

        tot_reward += reward
        steps += 1

        state, action = state_next, action_next

    episode_rewards[ep] = tot_reward
    steps_to_completion[ep] = steps

    if (ep+1)%print_freq == 0 and plot_heat:
        clear_output(wait=True)
        plot_Q(Q, message = "Episode %d: Reward: %f, Steps: %.2f,
Qmax: %.2f, Qmin: %.2f"%(ep+1, np.mean(episode_rewards[ep-
print_freq+1:ep]),
np.mean(steps_to_completion[ep-print_freq+1:ep]),
Q.max(), Q.min()))

    return Q, episode_rewards, steps_to_completion

```

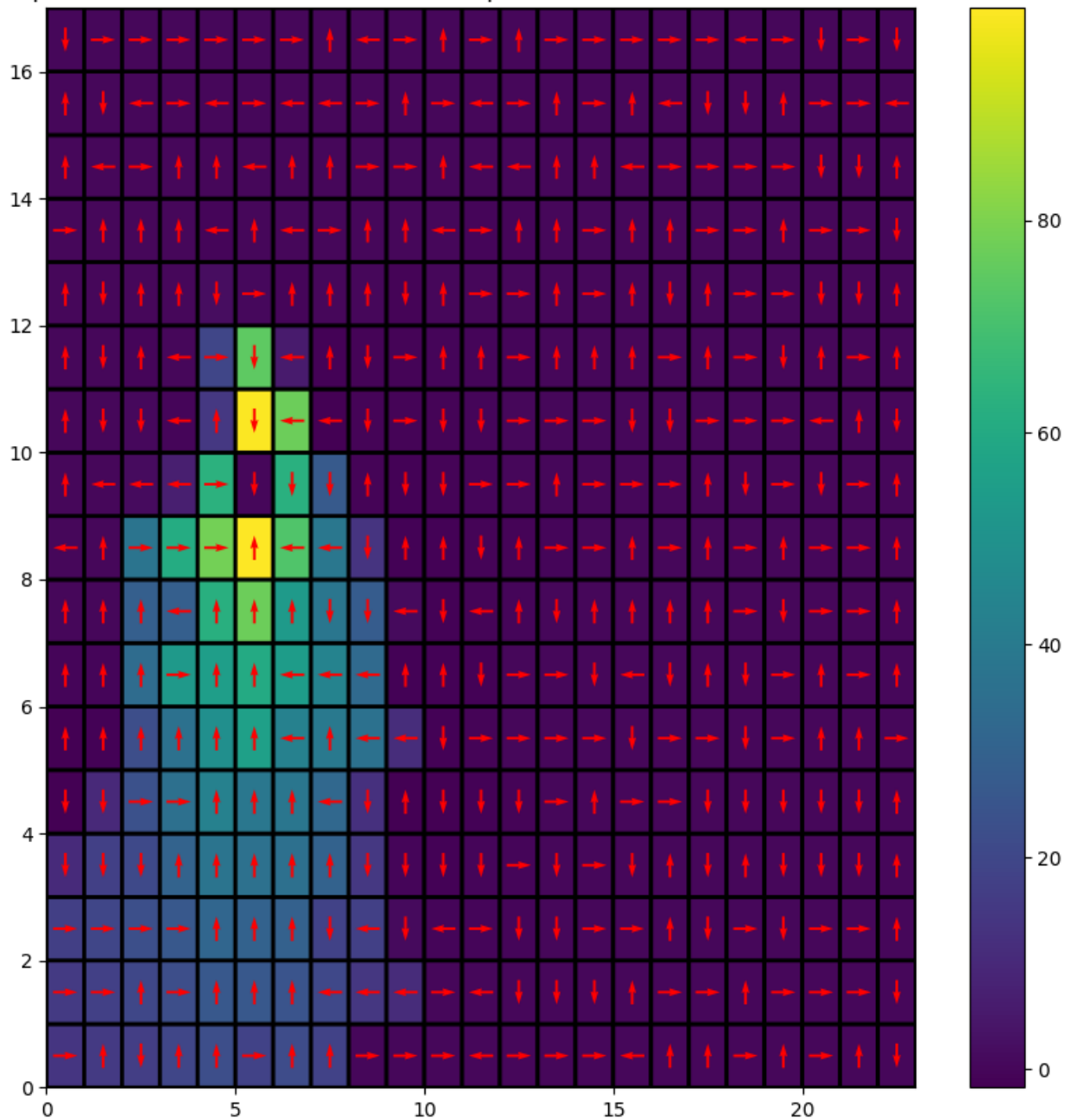
Q Learning with ϵ Greedy Policy

```

Q, rewards, steps = qlearning(env, Q, gamma = gamma, plot_heat=True,
choose_action= choose_action_epsilon)

```


Episode 10000: Reward: 23.575758, Steps: 7278.49, Qmax: 99.97, Qmin: -2.12



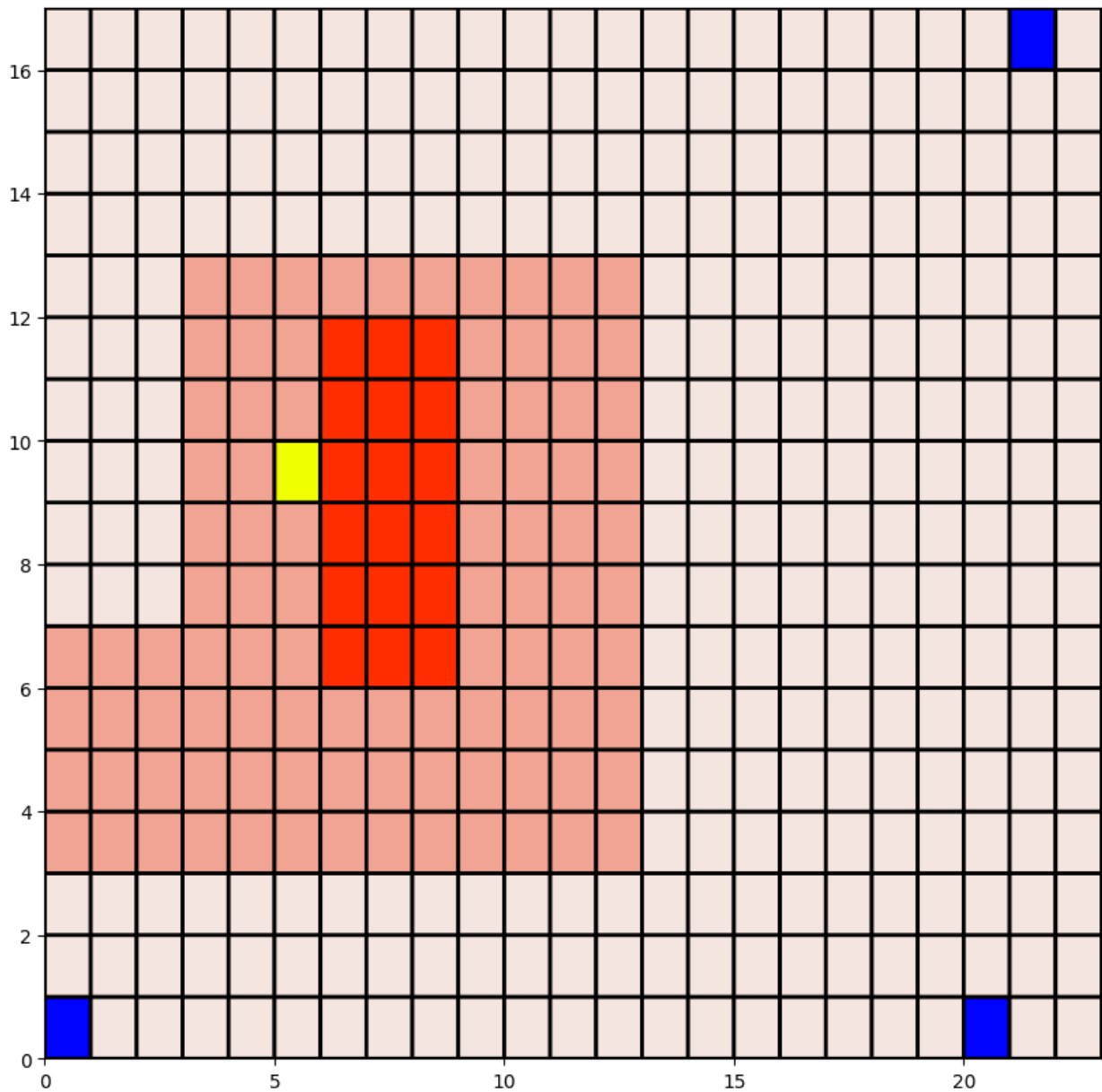
100%|██████████| 10000/10000 [11:50<00:00, 14.07it/s]

Visualizing the policy

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

```
from time import sleep  
state = env.reset()
```

```
done = False
steps = 0
tot_reward = 0
while not done:
    clear_output(wait=True)
    state, reward, done = env.step(Q[state[0], state[1]].argmax())
    plt.figure(figsize=(10, 10))
    env.render(ax=plt, render_agent=True)
    plt.show()
    steps += 1
    tot_reward += reward
    sleep(0.2)
print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```



Steps: 14, Total Reward: 93

Q Learning with Softmax Policy

Initializing Q values and Hyperparameters again

```
# initialize Q-value
Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
len(env.action_space)))

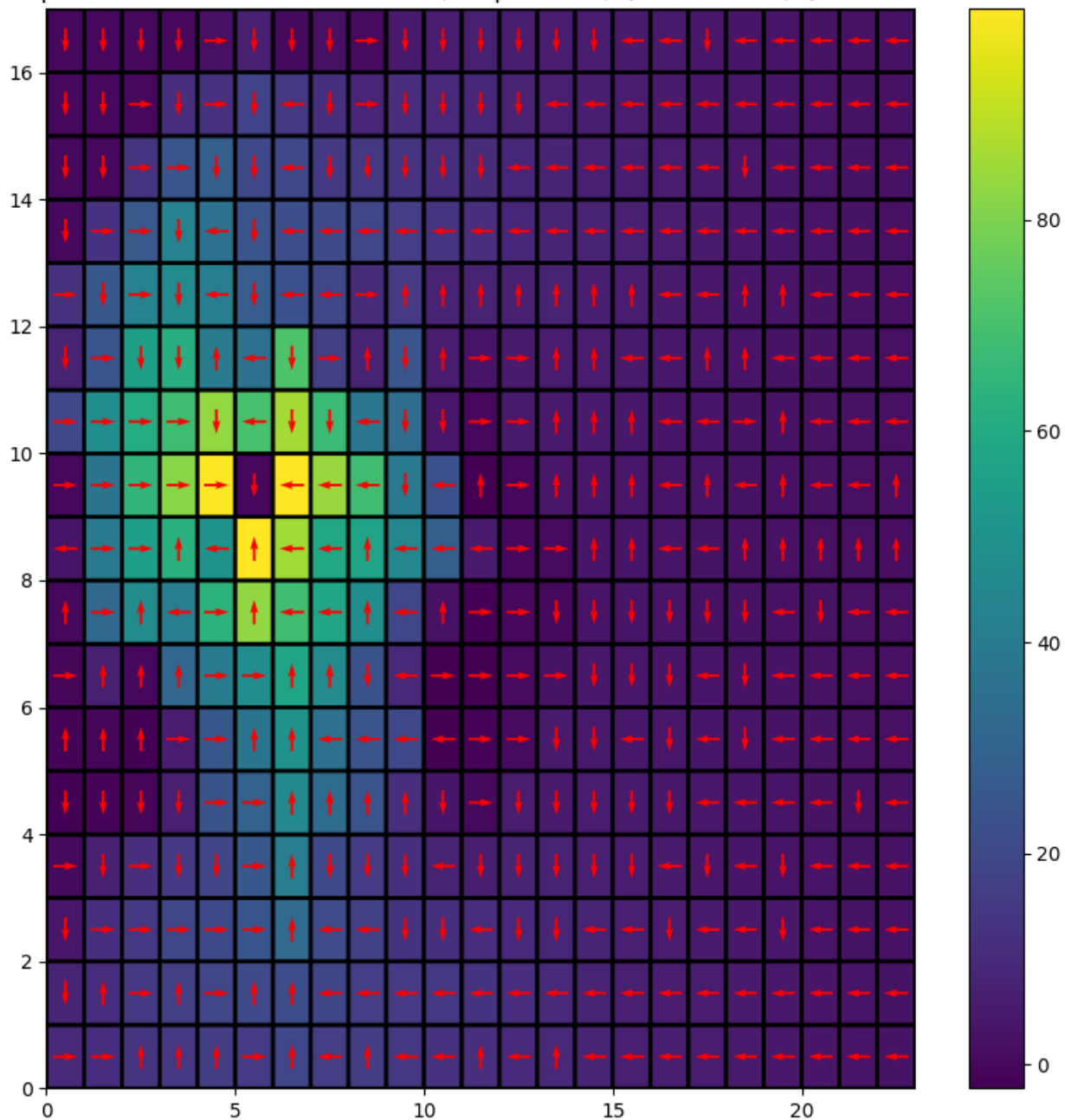
alpha0 = 0.4
gamma = 0.9
```

```
episodes = 10000
```

```
epsilon0 = 0.1
```

```
Q, rewards, steps = qlearning(env, Q, gamma = gamma, plot_heat=True,  
choose_action= choose_action_softmax)
```

Episode 10000: Reward: 87.696970, Steps: 52.16, Qmax: 100.00, Qmin: -3.00



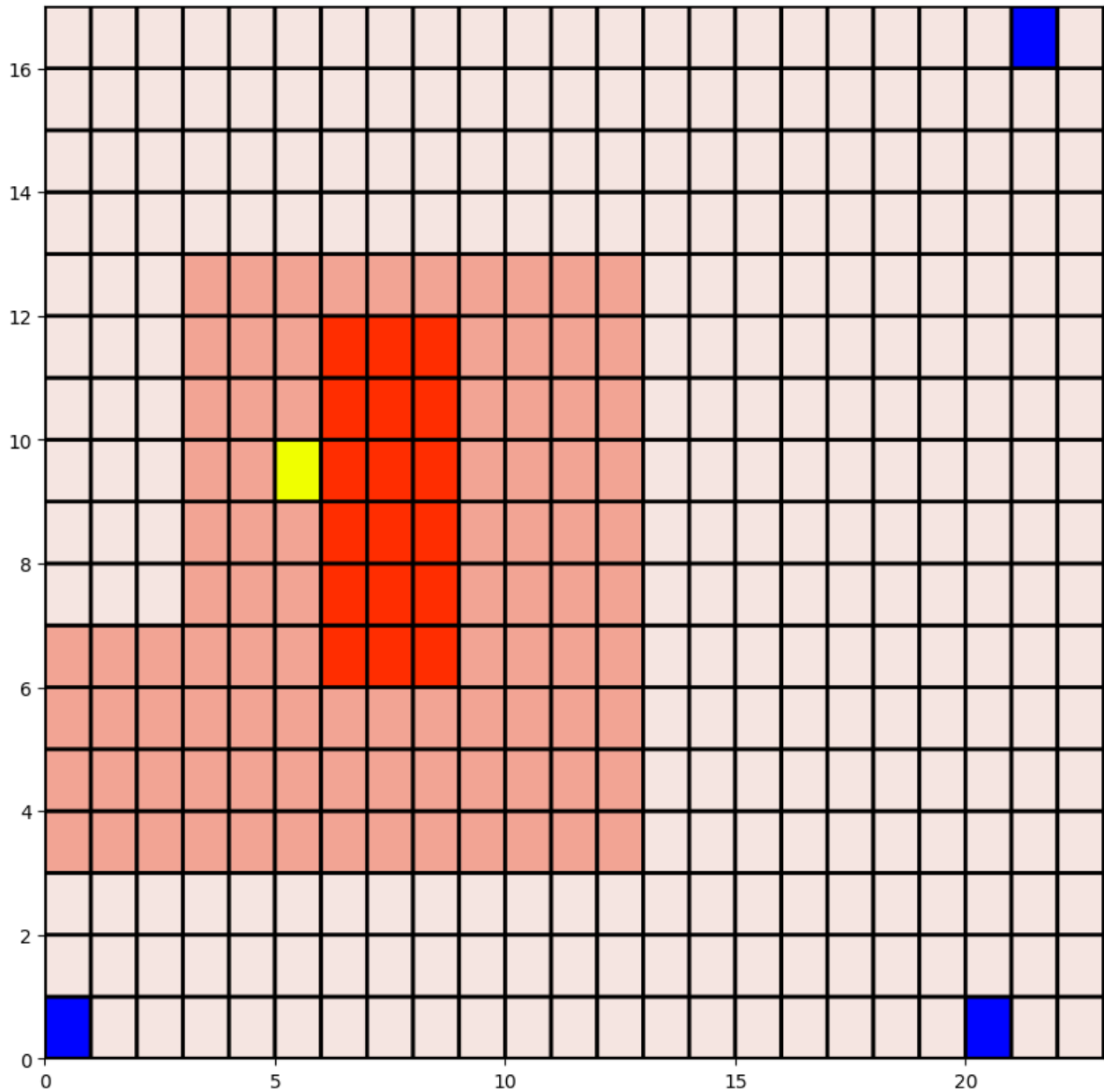
100%|██████████| 10000/10000 [00:46<00:00, 213.59it/s]

Visualizing the policy

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

```
from time import sleep

state = env.reset()
done = False
steps = 0
tot_reward = 0
while not done:
    clear_output(wait=True)
    state, reward, done = env.step(Q[state[0], state[1]].argmax())
    plt.figure(figsize=(10, 10))
    env.render(ax=plt, render_agent=True)
    plt.show()
    steps += 1
    tot_reward += reward
    sleep(0.2)
print("Steps: %d, Total Reward: %d"%(steps, tot_reward))
```



Steps: 16, Total Reward: 91

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.

```

num_expts = 5
reward_avgs, steps_avgs = [], []

for i in range(num_expts):
    print("Experiment: %d"%(i+1))
    Q = np.zeros((env.grid.shape[0], env.grid.shape[1],
len(env.action_space)))
    rg = np.random.RandomState(i)

    # TODO: run qlearning, store metrics
    Q, rewards, steps = qlearning(env, Q, gamma = gamma,
plot_heat=False, choose_action = choose_action_softmax)
    reward_avgs.append(rewards)
    steps_avgs.append(steps)

reward_avgs = np.array(reward_avgs)
steps_avgs = np.array(steps_avgs)
reward_avgs = np.mean(reward_avgs, axis=0)
steps_avgs = np.mean(steps_avgs, axis=0)

Experiment: 1
100%|██████████| 10000/10000 [00:23<00:00, 421.36it/s]

Experiment: 2
100%|██████████| 10000/10000 [00:36<00:00, 276.96it/s]

Experiment: 3
100%|██████████| 10000/10000 [00:22<00:00, 435.52it/s]

Experiment: 4
100%|██████████| 10000/10000 [00:24<00:00, 402.52it/s]

Experiment: 5
100%|██████████| 10000/10000 [00:26<00:00, 380.49it/s]

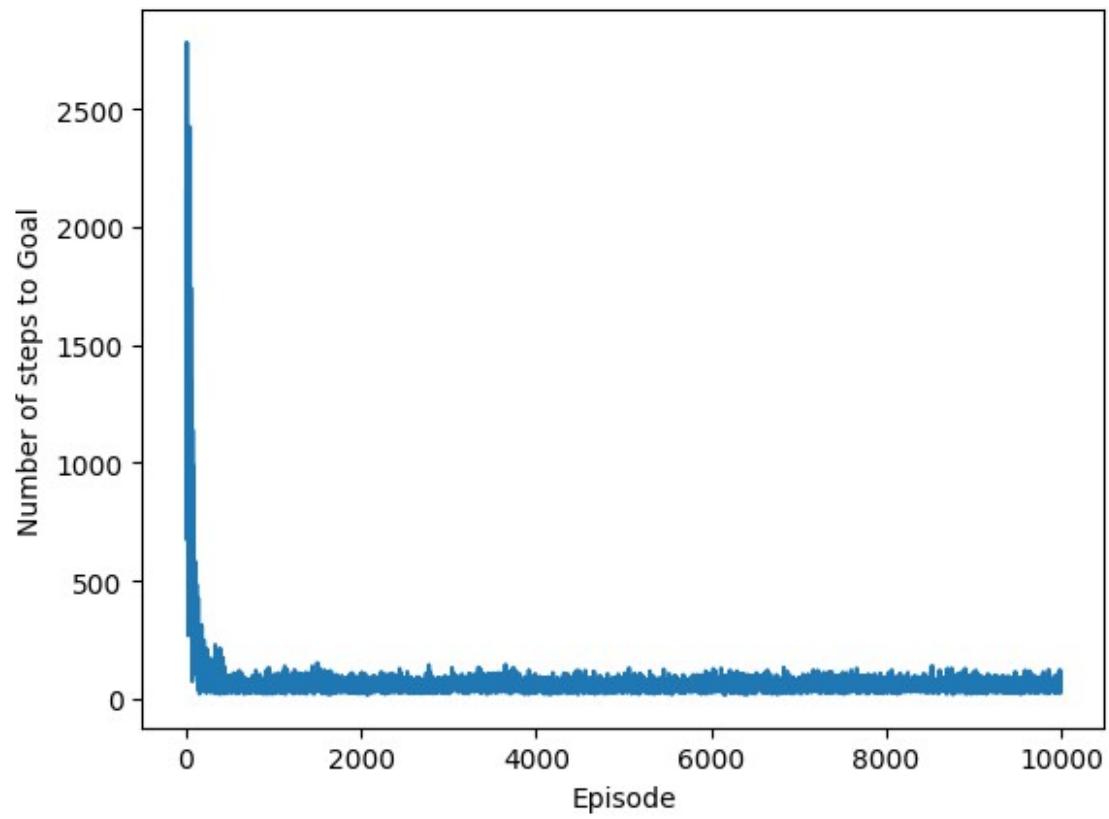
# TODO: visualize individual metrics vs episode count (averaged across
multiple run(s))

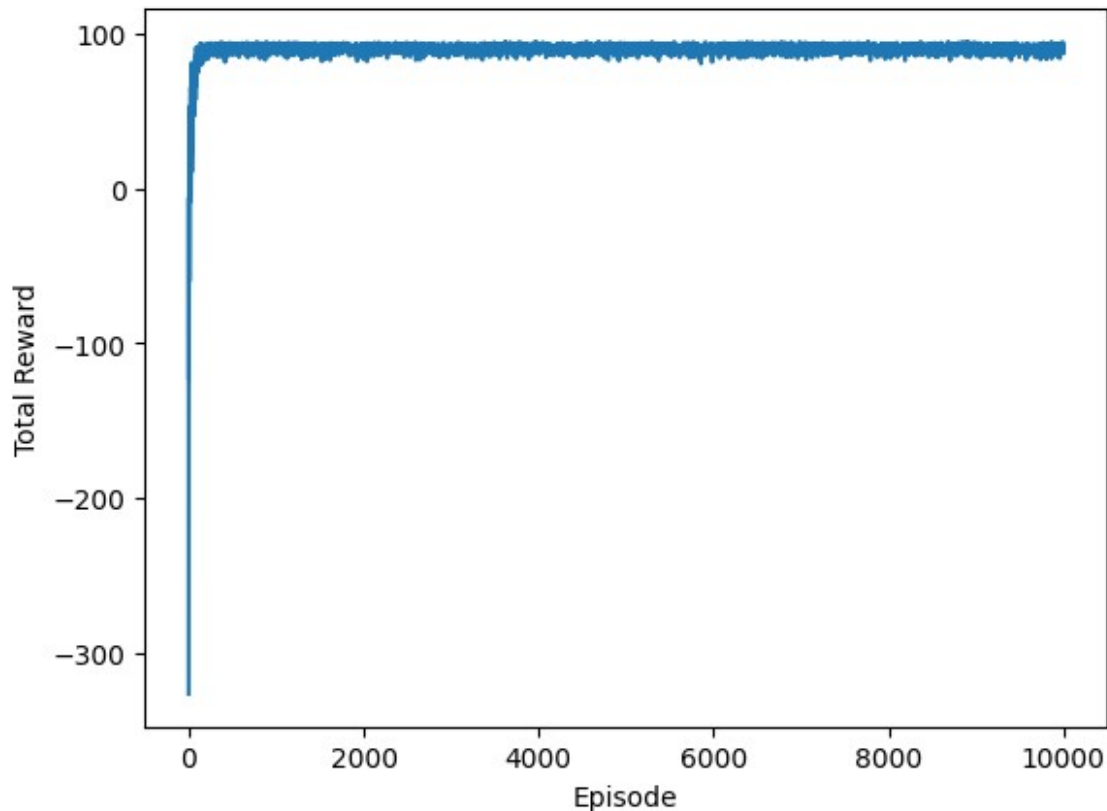
plt.figure()
plt.plot(steps_avgs)
plt.xlabel('Episode')
plt.ylabel('Number of steps to Goal')
plt.show()

plt.figure()
plt.plot(reward_avgs)
plt.xlabel('Episode')

```

```
plt.ylabel('Total Reward')  
plt.show()
```





TODO: What differences do you observe between the policies learnt by Q Learning and SARSA (if any).

- From the visual observation of the actions taken by the agent, we can see that the agent trained using SARSA algorithm takes a safer route when compared to the agent trained using Q Learning algorithm. A safer route here means a path that passes through the states that allow the agent to incur lesser negative rewards.
- Both algorithms converge to their respective optimal policies and perform equally good, which is to minimize the number of steps taken to reach the end goal. A subtle observation from the average rewards and steps graph is that the Q learning algorithm has lesser variation from the optimal number of steps as compared to the SARSA algorithm.
- We can also observe both the algorithms converge to their respective optimal policies at almost the same rate.
- In both SARSA and Q learning algorithms, the softmax policy performs better and runs faster compared to the ϵ greedy policy.
- Also we can infer from theory that, SARSA is an on-policy algorithm whereas Q learning is off-policy.

```
!pip install nbconvert
!sudo apt-get install texlive-xetex texlive-fonts-recommended texlive-plain-generic
```

Requirement already satisfied: nbconvert in c:\python311\lib\site-packages (7.2.8)

Requirement already satisfied: beautifulsoup4 in c:\python311\lib\site-packages (from nbconvert) (4.8.0)

Requirement already satisfied: bleach in c:\python311\lib\site-packages (from nbconvert) (5.0.1)

Requirement already satisfied: defusedxml in c:\python311\lib\site-packages (from nbconvert) (0.7.1)

Requirement already satisfied: jinja2>=3.0 in c:\python311\lib\site-packages (from nbconvert) (3.1.2)

Requirement already satisfied: jupyter-core>=4.7 in c:\users\srika\appdata\roaming\python\python311\site-packages (from nbconvert) (5.1.0)

Requirement already satisfied: jupyterlab-pygments in c:\python311\lib\site-packages (from nbconvert) (0.2.2)

Requirement already satisfied: markupsafe>=2.0 in c:\python311\lib\site-packages (from nbconvert) (2.1.2)

Requirement already satisfied: mistune<3,>=2.0.3 in c:\python311\lib\site-packages (from nbconvert) (2.0.4)

Requirement already satisfied: nbclient>=0.5.0 in c:\python311\lib\site-packages (from nbconvert) (0.7.2)

Requirement already satisfied: nbformat>=5.1 in c:\python311\lib\site-packages (from nbconvert) (5.7.3)

Requirement already satisfied: packaging in c:\users\srika\appdata\roaming\python\python311\site-packages (from nbconvert) (22.0)

Requirement already satisfied: pandocfilters>=1.4.1 in c:\python311\lib\site-packages (from nbconvert) (1.5.0)

Requirement already satisfied: pygments>=2.4.1 in c:\users\srika\appdata\roaming\python\python311\site-packages (from nbconvert) (2.13.0)

Requirement already satisfied: tinycss2 in c:\python311\lib\site-packages (from nbconvert) (1.2.1)

Requirement already satisfied: traitlets>=5.0 in c:\users\srika\appdata\roaming\python\python311\site-packages (from nbconvert) (5.7.1)

Requirement already satisfied: platformdirs>=2.5 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-core>=4.7->nbconvert) (2.6.0)

Requirement already satisfied: pywin32>=1.0 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-core>=4.7->nbconvert) (305)

Requirement already satisfied: jupyter-client>=6.1.12 in c:\users\srika\appdata\roaming\python\python311\site-packages (from nbclient>=0.5.0->nbconvert) (7.4.8)

Requirement already satisfied: fastjsonschema in c:\python311\lib\site-packages (from nbformat>=5.1->nbconvert) (2.16.2)

Requirement already satisfied: jsonschema>=2.6 in c:\python311\lib\site-packages (from nbformat>=5.1->nbconvert) (4.17.3)

Requirement already satisfied: soupsieve>=1.2 in c:\python311\lib\site-packages (from beautifulsoup4->nbconvert) (2.3.2.post1)

Requirement already satisfied: six>=1.9.0 in c:\users\srika\appdata\roaming\python\python311\site-packages (from bleach->nbconvert) (1.16.0)
Requirement already satisfied: webencodings in c:\python311\lib\site-packages (from bleach->nbconvert) (0.5.1)
Requirement already satisfied: attrs>=17.4.0 in c:\python311\lib\site-packages (from jsonschema>=2.6->nbformat>=5.1->nbconvert) (22.2.0)
Requirement already satisfied: pyparsing!=0.17.0,!=0.17.1,!=0.17.2,>=0.14.0 in c:\python311\lib\site-packages (from jsonschema>=2.6->nbformat>=5.1->nbconvert) (0.19.3)
Requirement already satisfied: entrypoints in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-client>=6.1.12->nbclient>=0.5.0->nbconvert) (0.4)
Requirement already satisfied: nest-asyncio>=1.5.4 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-client>=6.1.12->nbclient>=0.5.0->nbconvert) (1.5.6)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-client>=6.1.12->nbclient>=0.5.0->nbconvert) (2.8.2)
Requirement already satisfied: pyzmq>=23.0 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-client>=6.1.12->nbclient>=0.5.0->nbconvert) (24.0.1)
Requirement already satisfied: tornado>=6.2 in c:\users\srika\appdata\roaming\python\python311\site-packages (from jupyter-client>=6.1.12->nbclient>=0.5.0->nbconvert) (6.2)

[notice] A new release of pip is available: 23.2.1 -> 24.0
[notice] To update, run: python.exe -m pip install --upgrade pip
ERROR: Could not find a version that satisfies the requirement
texlive-xetex (from versions: none)
ERROR: No matching distribution found for texlive-xetex

[notice] A new release of pip is available: 23.2.1 -> 24.0
[notice] To update, run: python.exe -m pip install --upgrade pip

!jupyter nbconvert --to html "/content/drive/MyDrive/Colab
Notebooks/CS6700_Tutorial_4_QLearning_SARSA_ROLLNUMBER.ipynb"