# CS6700: Reinforcement Learning

# Q-learning Report

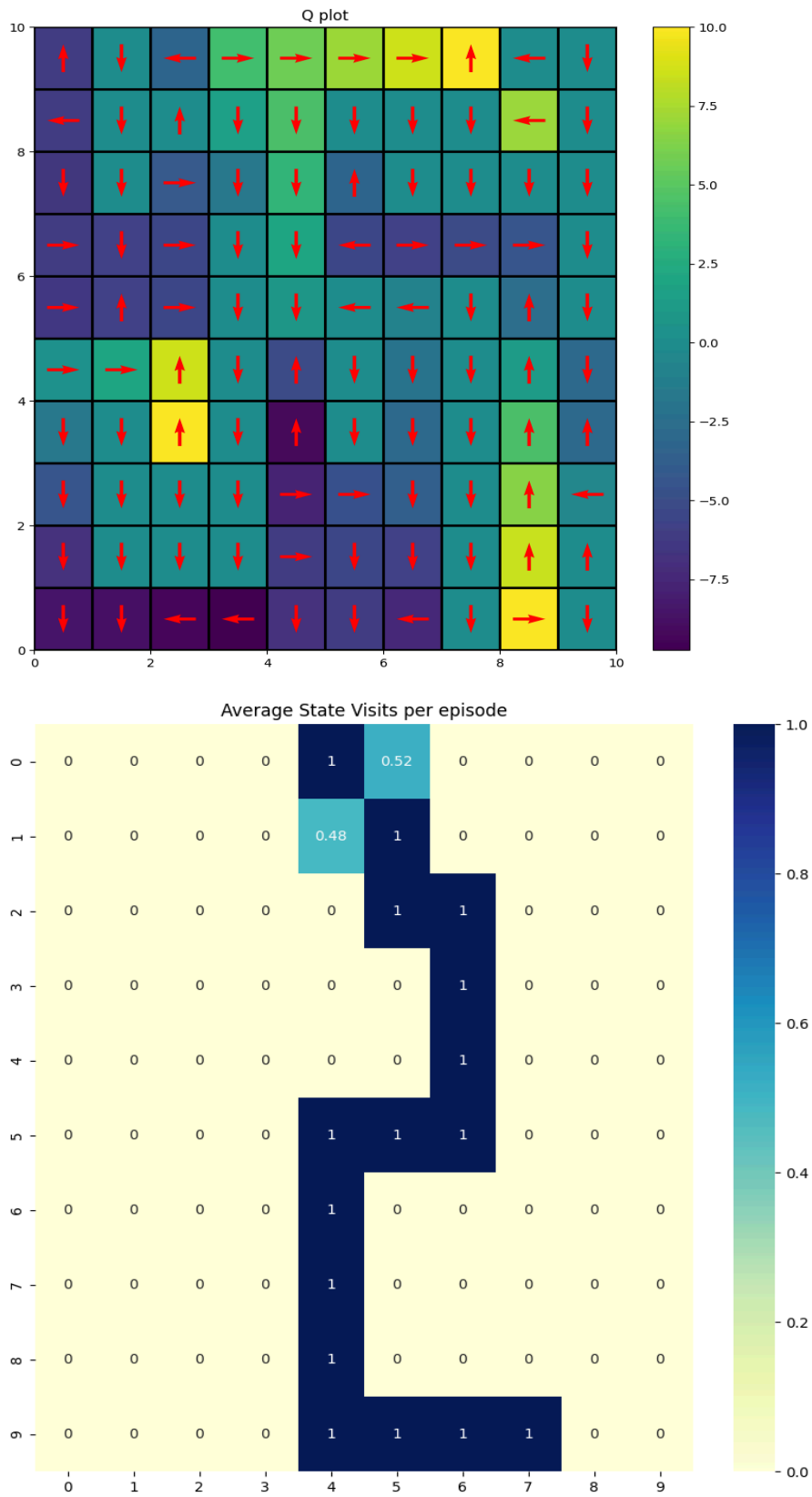- Shaunak Mujumdar ( CH21B062)
- Vedant Bhume (ED21B017)

Both of the policies with specific hyper parameters were trained for 1000 episodes. The best hyper parameters were chosen based on high mean reward of the policy for the last 100 episodes of training. The policy with the best hyper parameters was trained for 1000 episodes on 5 different experiments and the final results are shown.
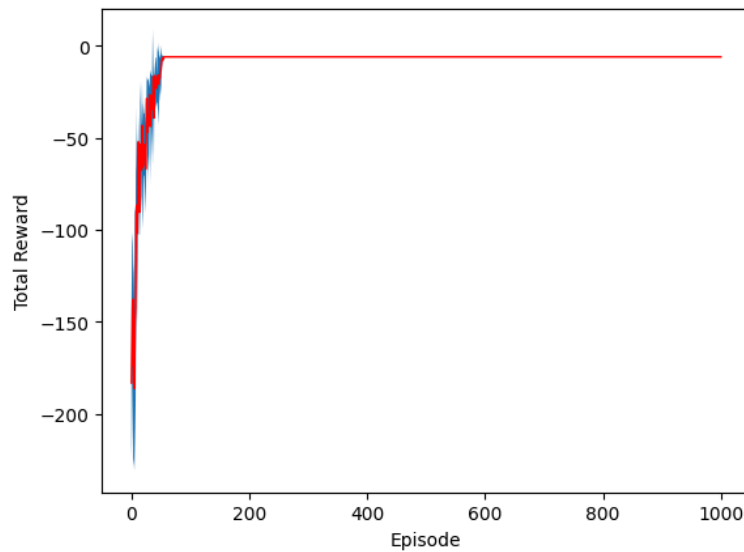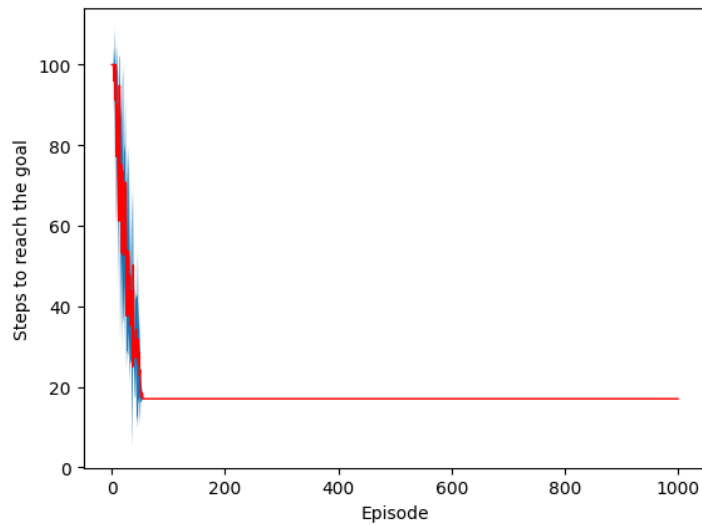
Because of the time constraints, it wasn't possible for us to explore a wide range of hyperparameters for all the given situations. We have only chosen a subset of hyperparameters and the algorithm will test out all possible combinations of these hyperparameters.

1. Alphas => [0.1,0.2,0.3,0.5,.7]
2. Gammas => [0.1,0.3,0.5,0.7,0.95]
3. Policy params => [0.01,0.05,0.3,0.5]

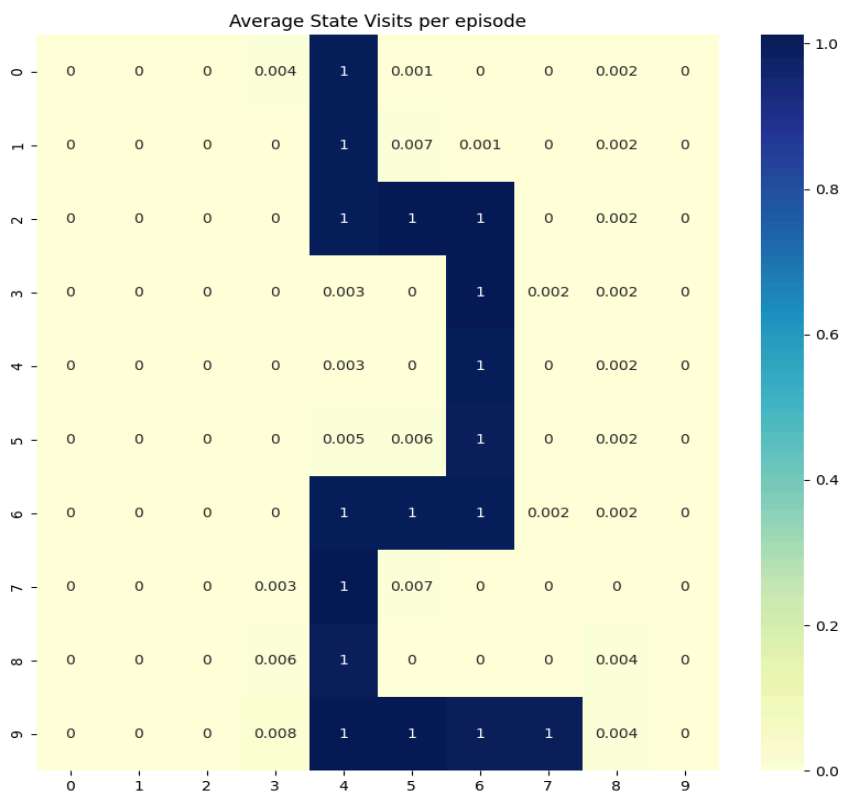| Start state | Wind | Transition Probability | Action Policy | Best Learning Rate ($\alpha$) | Best Discount Factor ($\gamma$) | Best Policy Parameter |
|---|---|---|---|---|---|---|
| [0,4] | False | 1 | Softmax | 0.7 | 0.95 | 0.05 |
| [0,4] | False | 1 | E-greedy | 0.7 | 0.95 | 0.01 |
| [3,6] | False | 1 | Softmax | 0.7 | 0.7 | 0.01 |
| [3,6] | False | 1 | E-greedy | 0.7 | 0.5 | 0.01 |
| [0,4] | False | 0.7 | Softmax | 0.5 | 0.95 | 0.01 |
| [0,4] | False | 0.7 | E-greedy | 0.3 | 0.95 | 0.01 |
| [3,6] | False | 0.7 | Softmax | 0.5 | 0.95 | 0.05 |
| [3,6] | False | 0.7 | E-greedy | 0.3 | 0.95 | 0.01 |
| [0,4] | True | 1 | Softmax | 0.7 | 0.95 | 0.01 |
| [0,4] | True | 1 | E-greedy | 0.7 | 0.95 | 0.01 |
| [3,6] | True | 1 | Softmax | 0.5 | 0.7 | 0.01 |
| [3,6] | True | 1 | E-greedy | 0.7 | 0.95 | 0.01 |

**1. Start state => [0,4] ; Wind => False ; P_transition = 1; Policy = Softmax:**
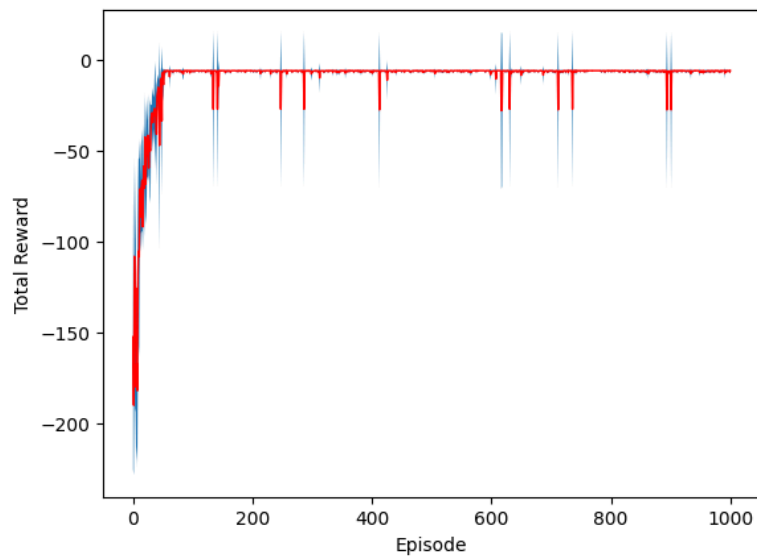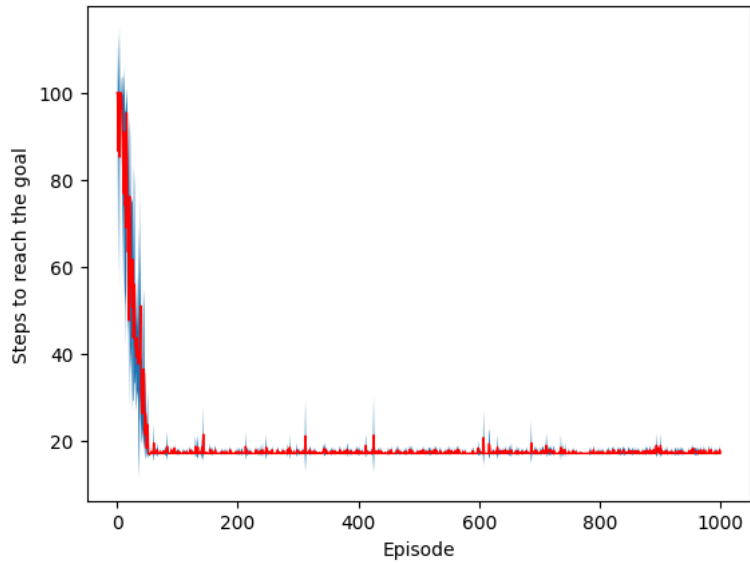
Q plot

Average State Visits per episode

By making p=1, we're removing all elements of stochasticity from the environment.The agent immediately learns a route to a goal state and keeps following this path for all following episodes. From the state visit heat map , you can see that policy never ends up accidentally exploring the environment and possibly finding the other goal states.
For all the experiments total reward and steps to reach the goal become constant after around 50 episodes and there is no standard deviation in the total reward and steps value.

## 2. Start state => [0,4] ; Wind => False ; P_transition = 1; Policy = E-greedy:



Q plot



Average State Visits per episode

From the state visited heatmap and reward/step vs episodes graphs, we can see that E-greedy approach shows far more random exploration than its softmax counterpart that was shown before this.

The mean number of steps is 17 and the mean reward is -4 (exactly like softmax).
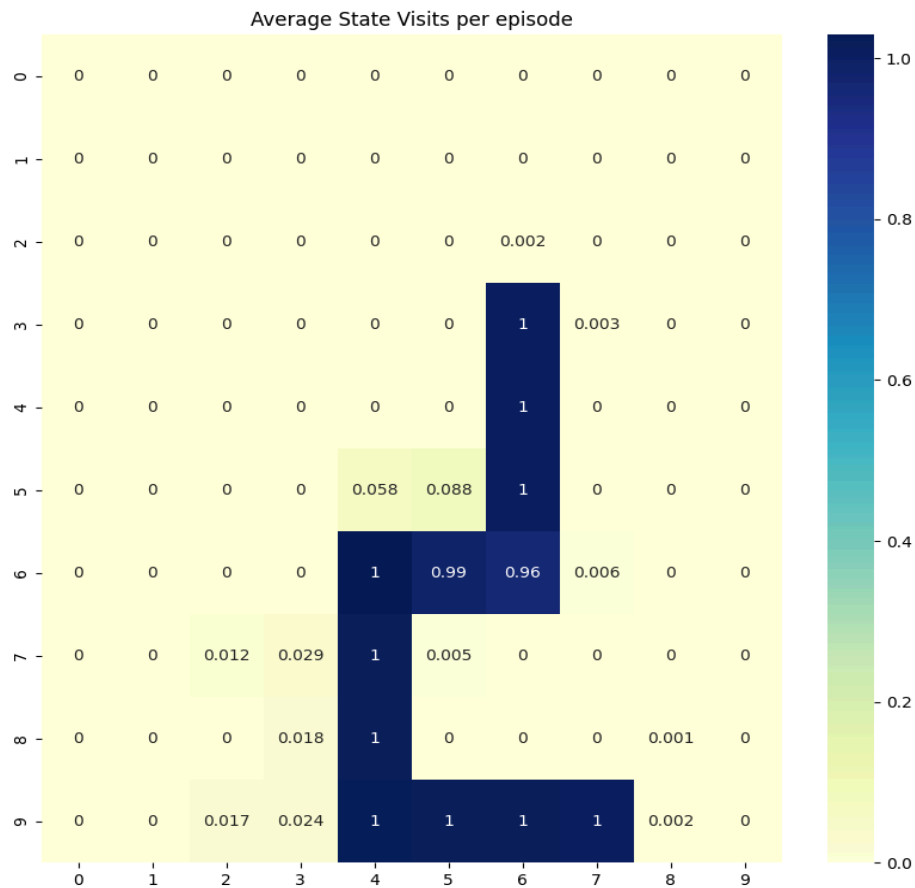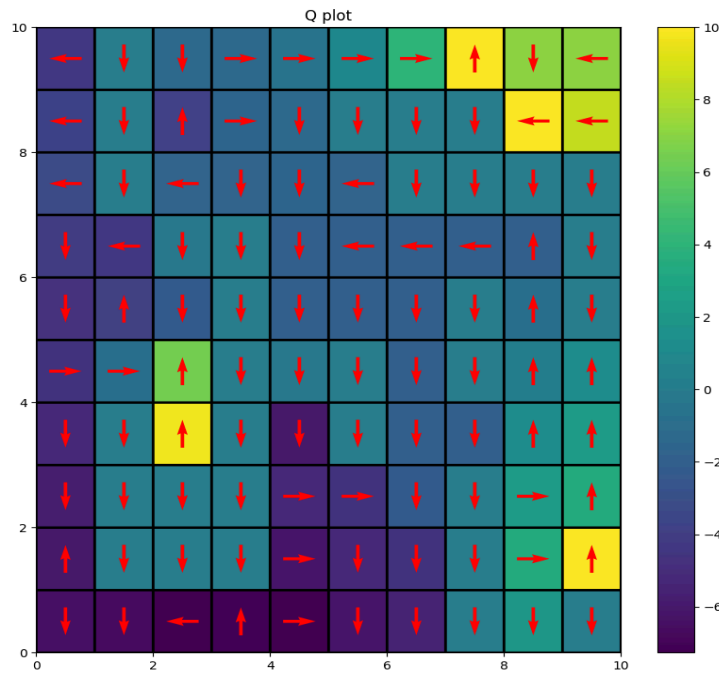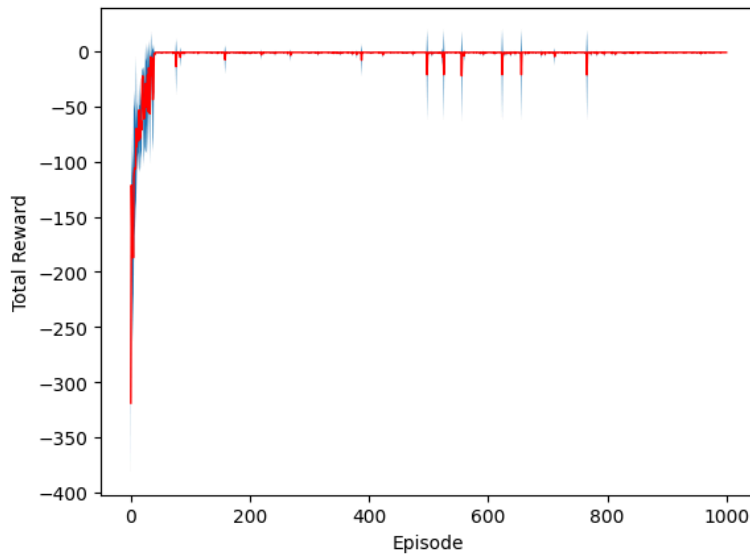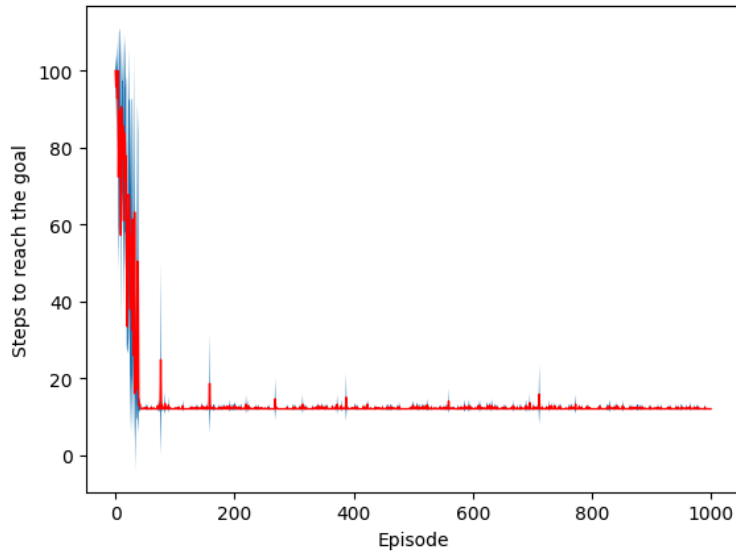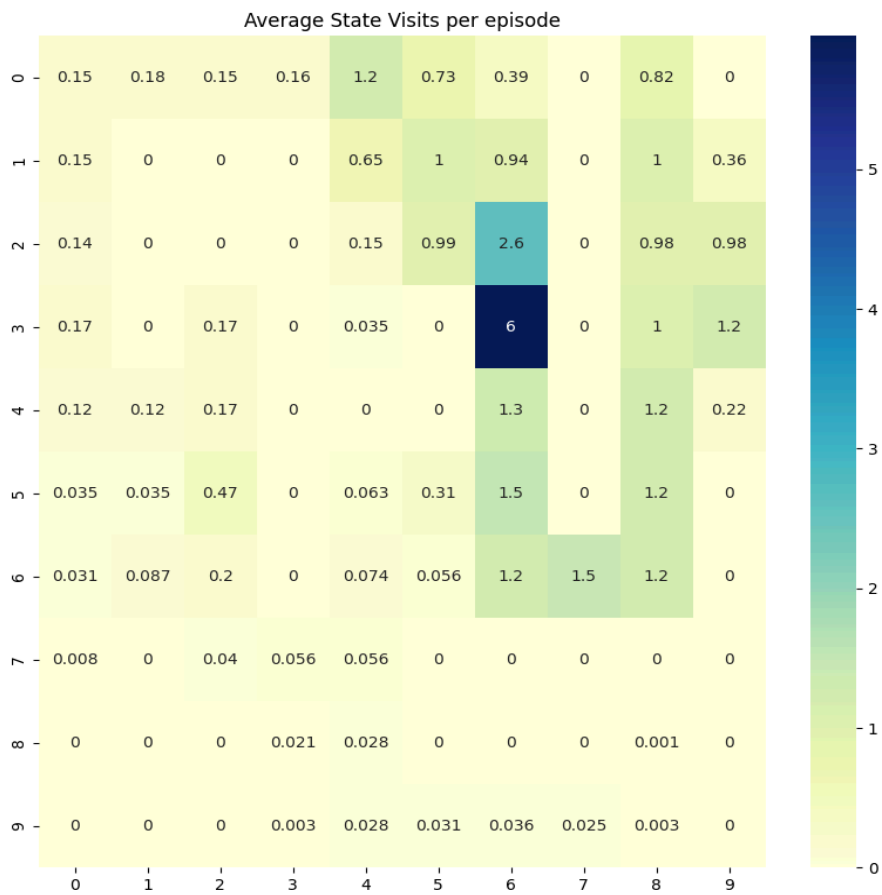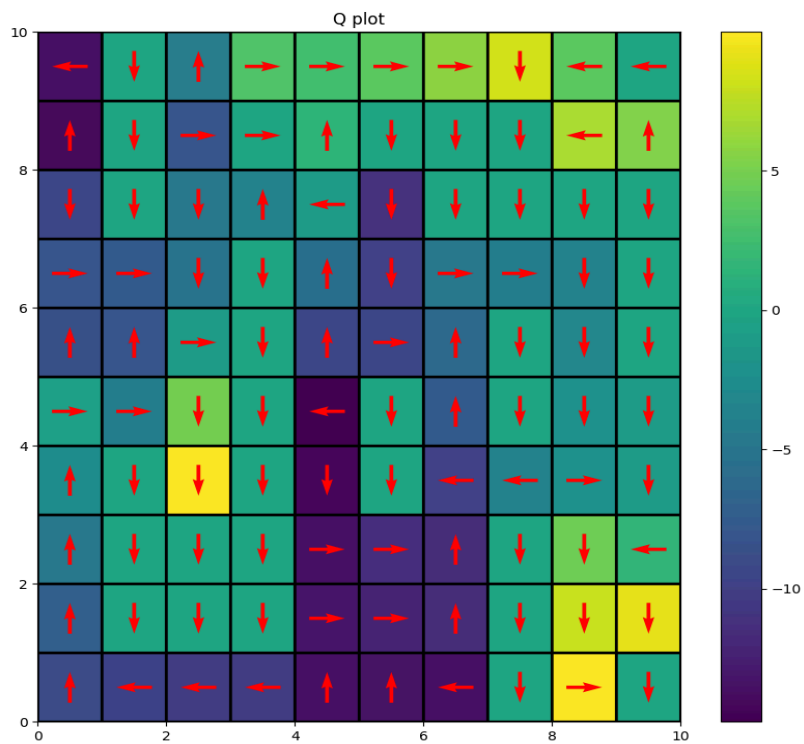
**3. Start state => [3,6] ; Wind => False ; P_transition = 1; Policy = Softmax :**



Q plot



Average State Visits per episode

The agent immediately learns a route to a goal state and keeps following this path for all following episodes. From the state visit heat map , you can see that policy never ends up accidentally exploring the environment and possibly finding the other goal states.

For all the experiments total reward and steps to reach the goal become constant after around 50 episodes and there is no standard deviation in the total reward and steps value. The mean number of steps is 10 and the mean reward is 0.

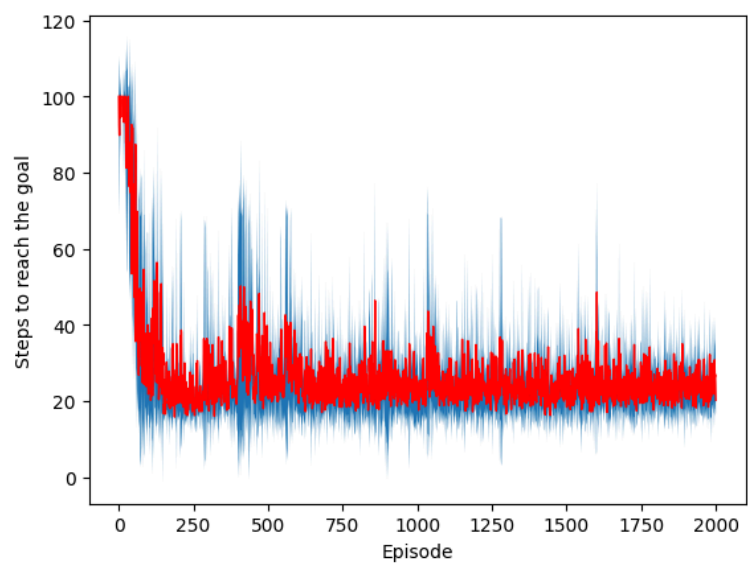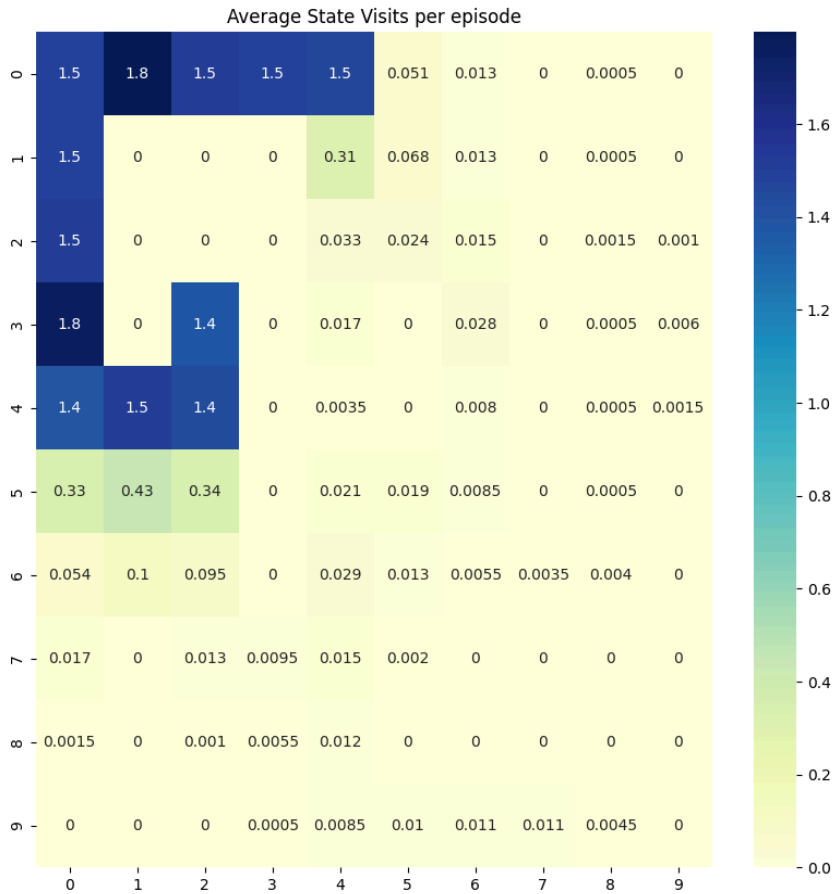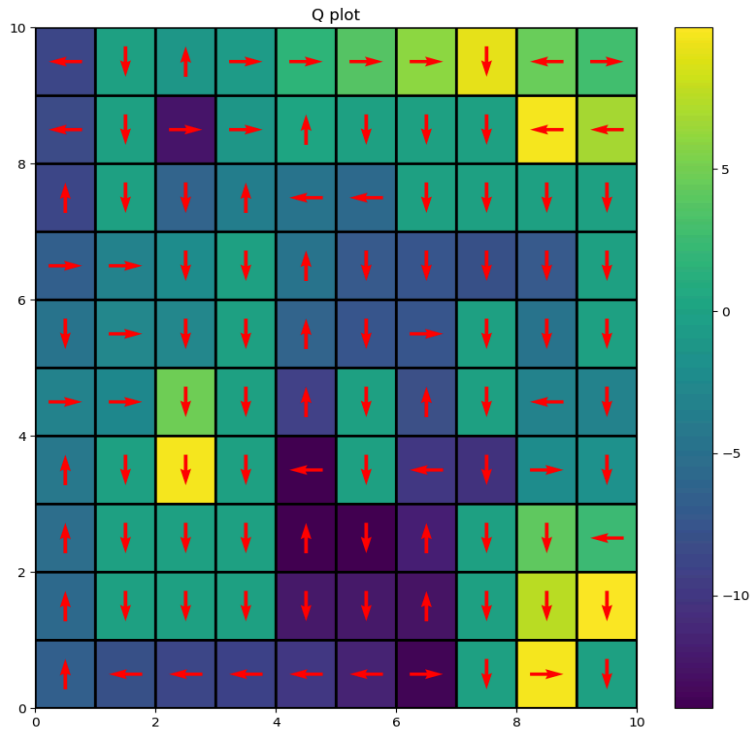**4. Start state => [3,6] ; Wind => False ; P_transition = 1; Policy = E-greedy:**



Q plot



Average State Visits per episode

Again just as in case with start case => [0,4],E-greedy approach shows far more random exploration than its softmax counterpart which is explicitly evident in the plots above. Instead of continuously exploiting the route that converges constantly, the E-greedy approach has the agent randomly explore the grid sometimes which accounts for the **random bumps present in the steps and rewards plots** and some **random non zero values in state visit heatmap.**
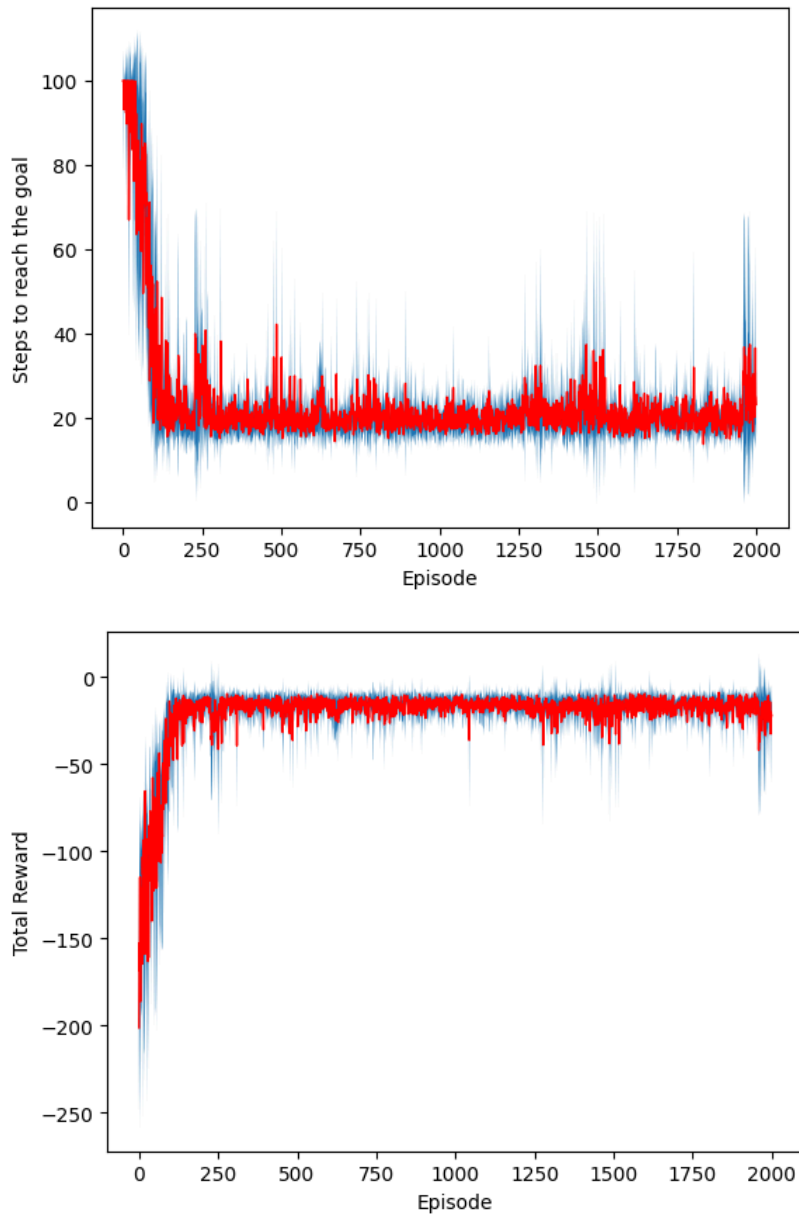
**5.  Start state => [0,4] ; Wind => False ; P_transition = 0.7; Policy =  Softmax :**



Q plot



Average State Visits per episode

**6. Start state => [0,4] ; Wind => False ; P_transition = 0.7; Policy = E-greedy:**


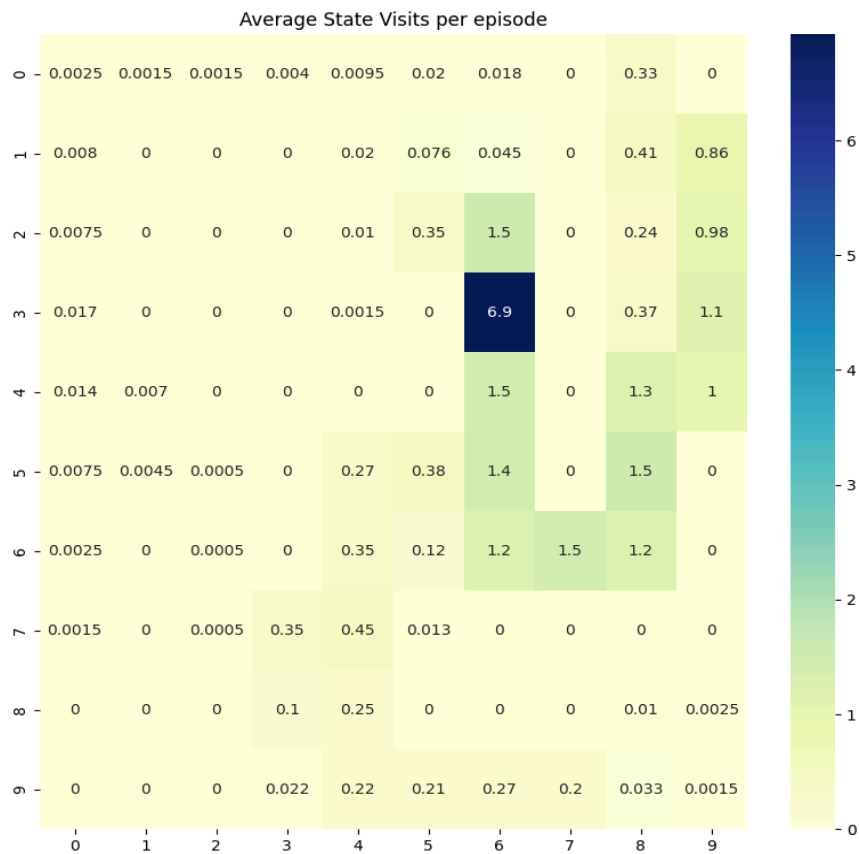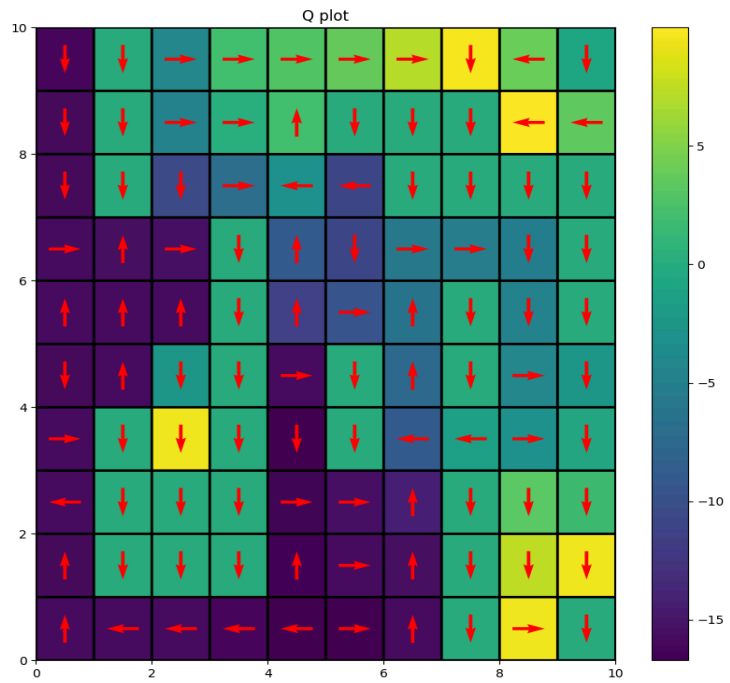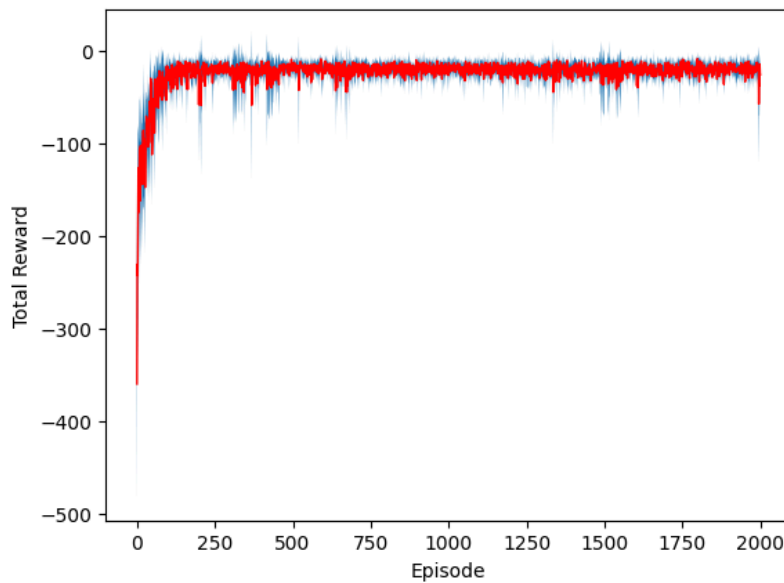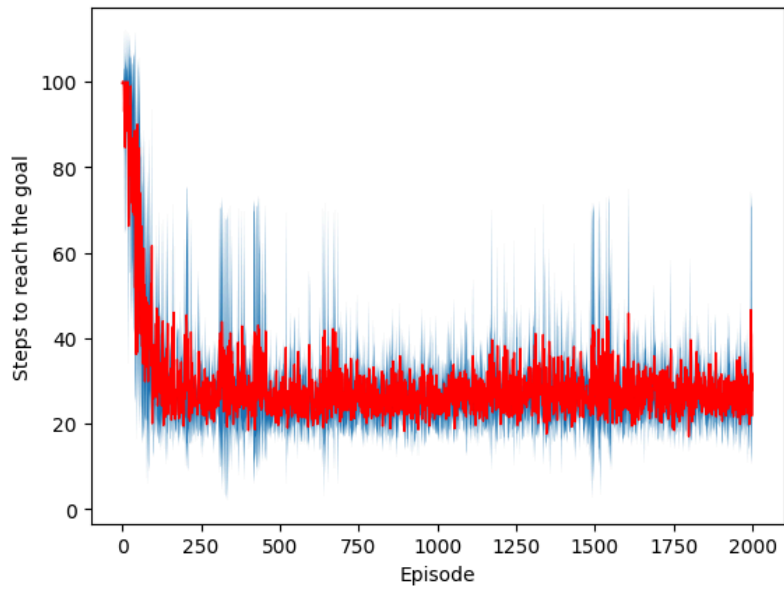
Q plot



Average State Visits per episode

From the state visit heat map, the agent learns the path to the nearest goal state. The path to the second nearest goal is also traced albeit very few times. This is in contrast to the SARSA algorithm where this second path is traced multiple more times. This indicates that the Q learning approach exploits the confirmed route to a goal state. The alternative route is taken primarily when the stochastic nature of the environment comes into play.

As the episodes increase, average reward converges to around -6 and average number of steps converges around 20.

**7. Start state => [3,6] ; Wind => False ; P_transition = 0.7; Policy = Softmax :**


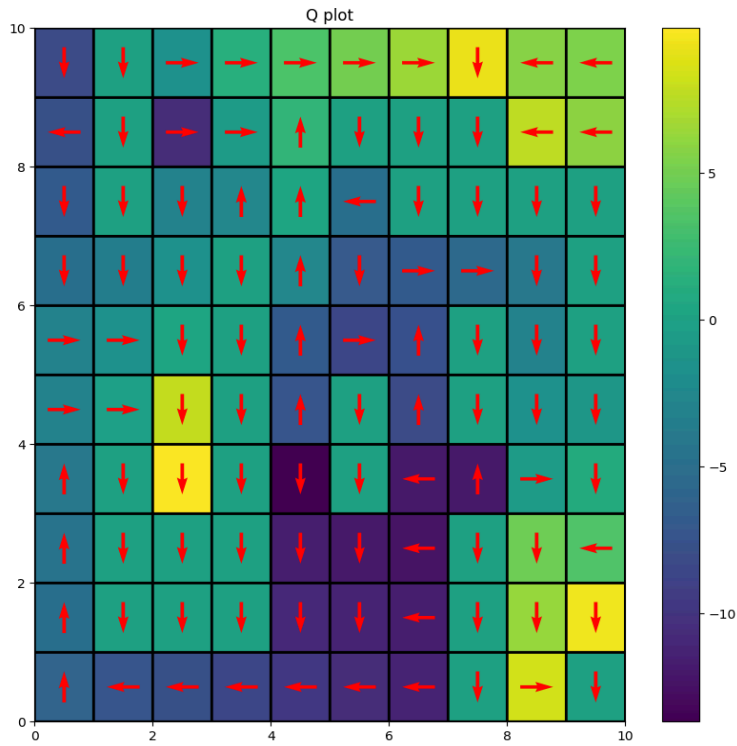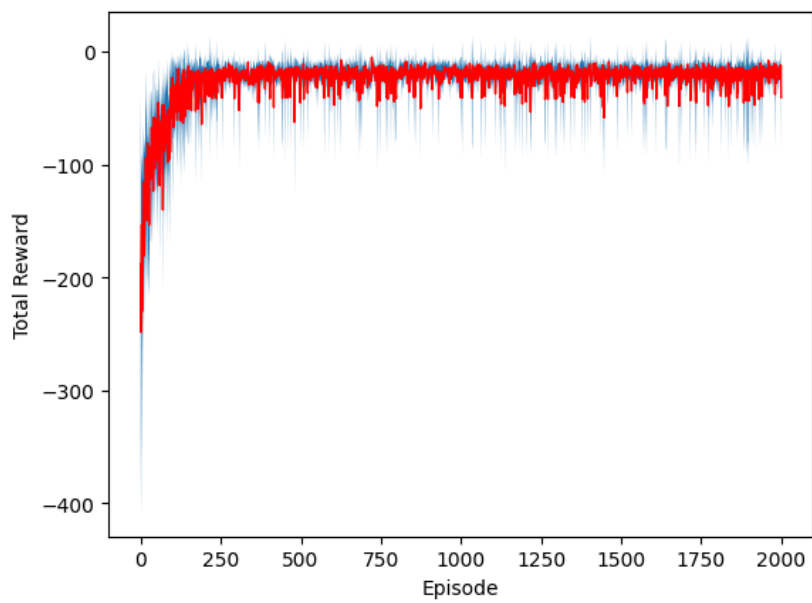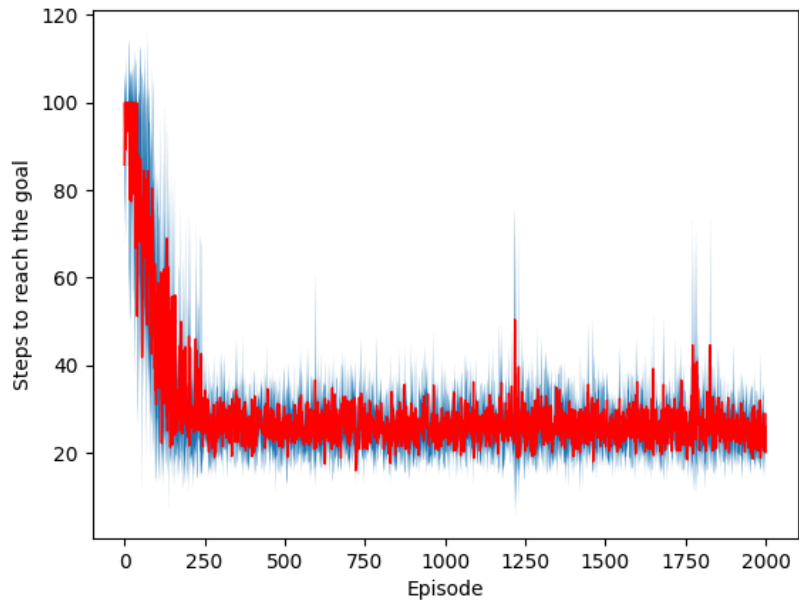
Q plot



Average State Visits per episode

This case is exactly similar to the first case, except the start state is located more in the middle of the environment.

As a result, the policy seems to equally prioritize between the two closest goal states unlike the first case.

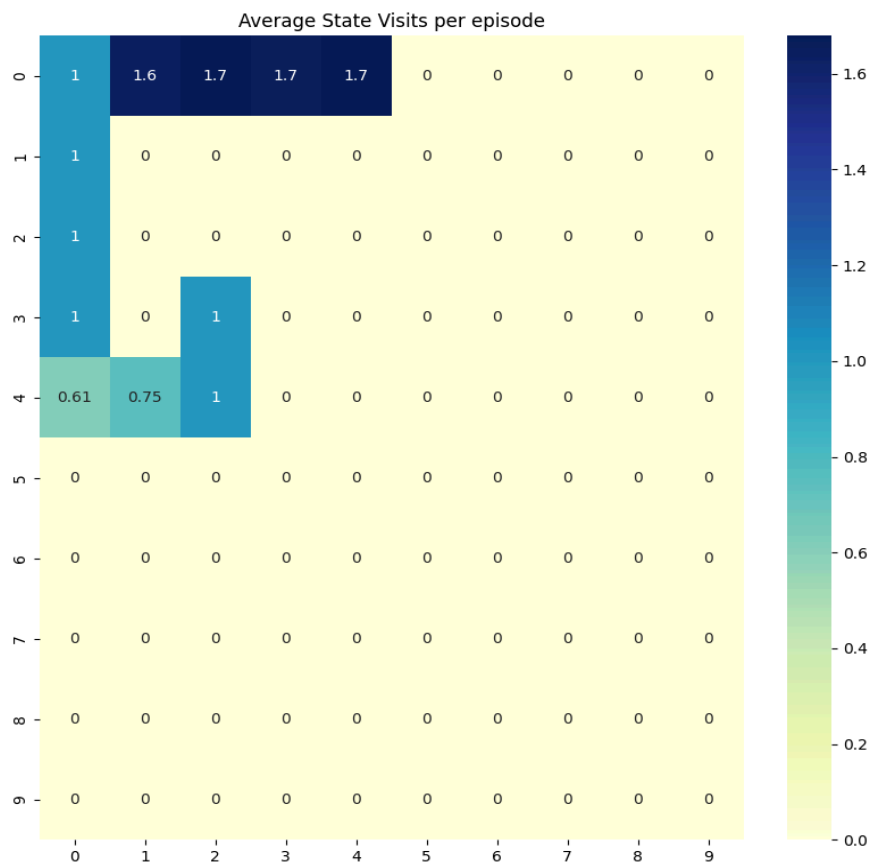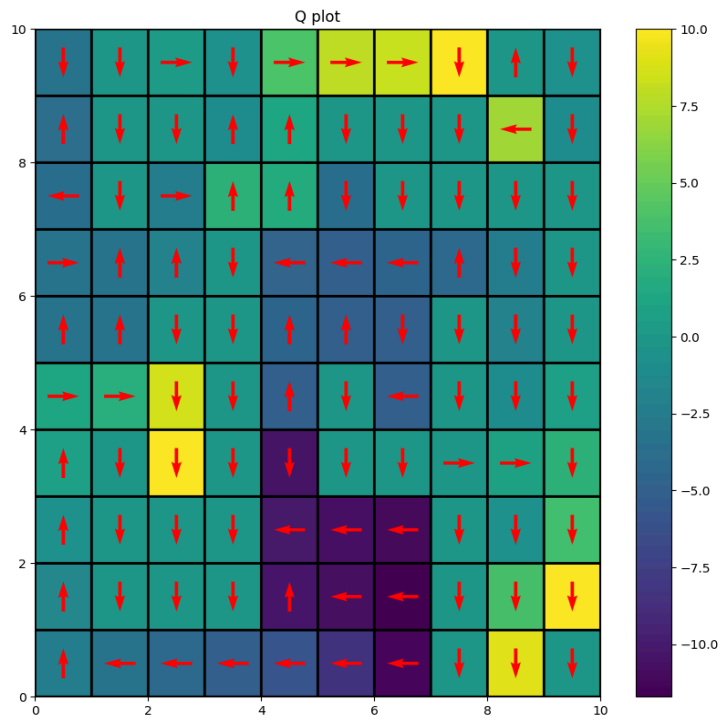The algorithm converges to 20 steps and -10 reward.

**8. Start state => [3,6] ; Wind => False ; P_transition = 0.7; Policy =  E-greedy :**
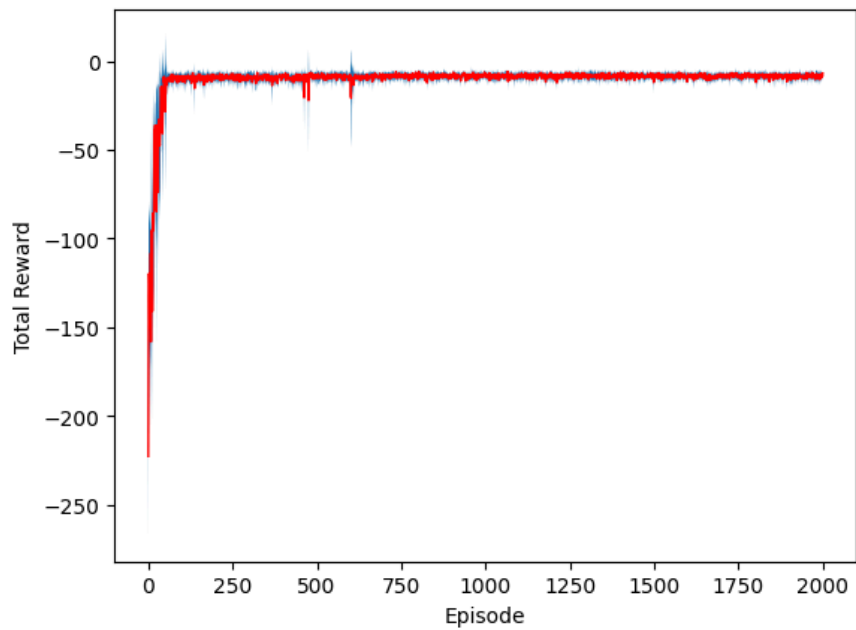


Q plot
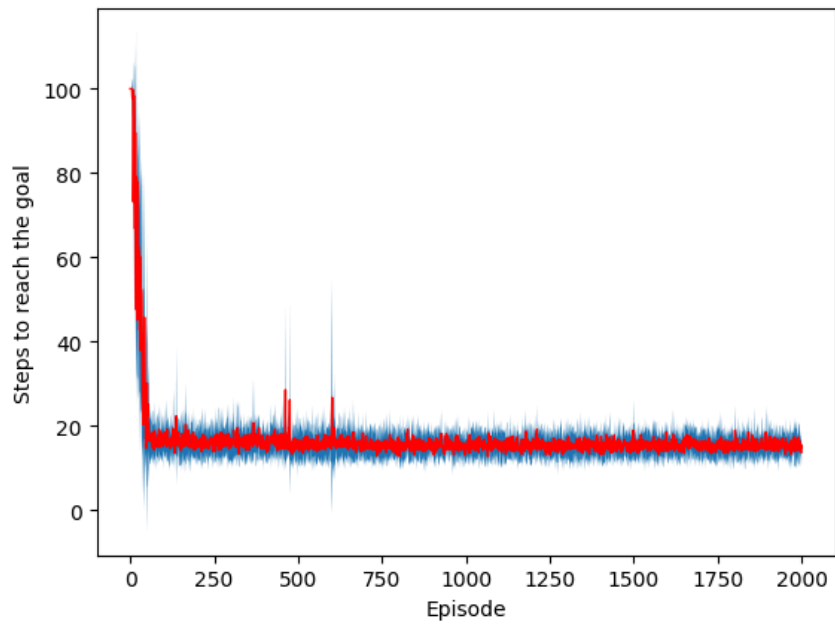


Average State Visits per episode

Large noise in the mean total reward and steps to reach the goal even with increasing episode number is because of the random exploration by E-greedy policy.

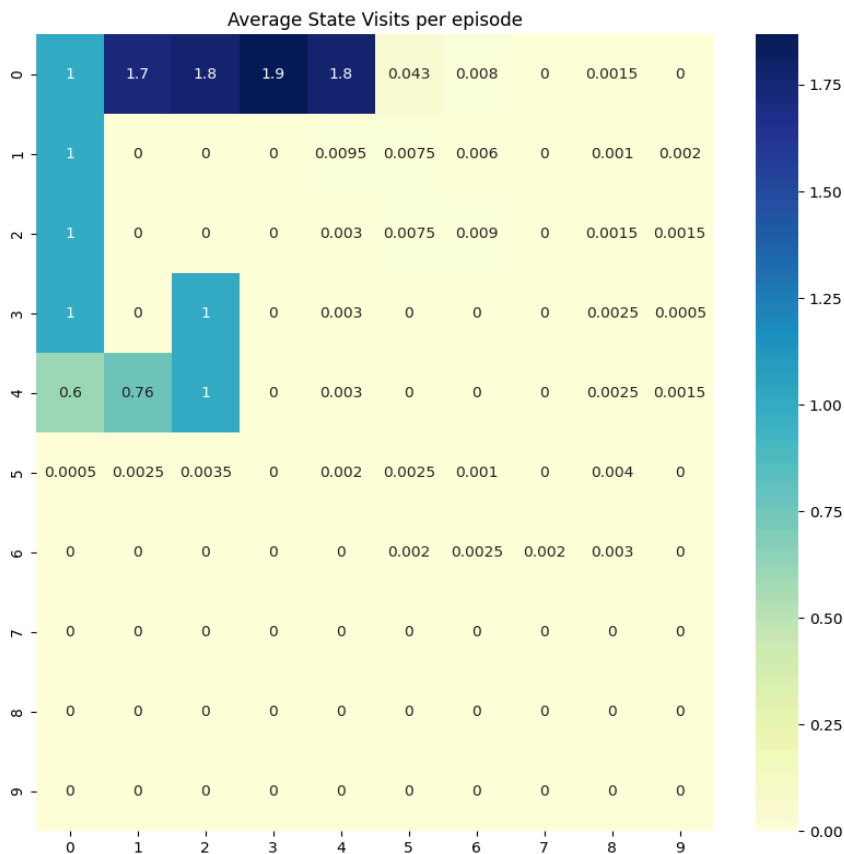**9. Start state => [0,4] ; Wind => True ; P_transition = 1 ; Policy = Softmax :**



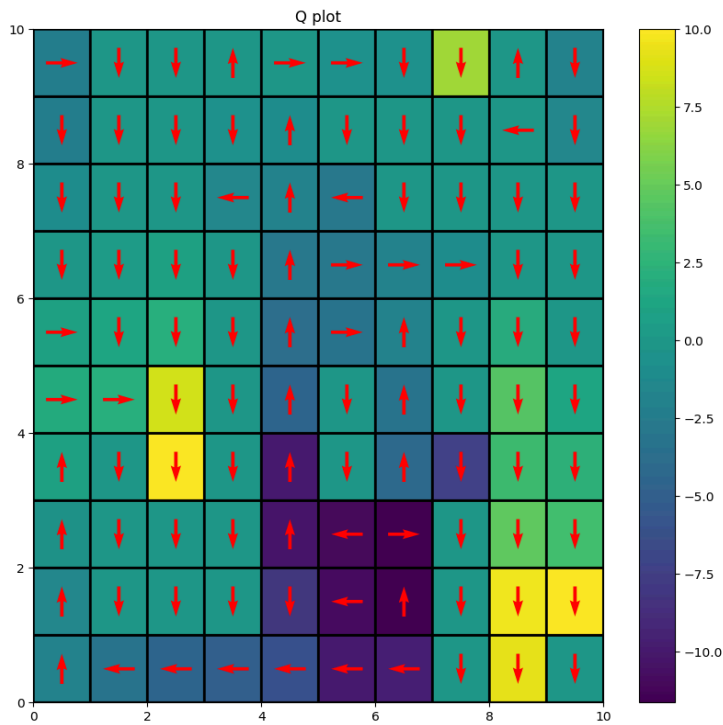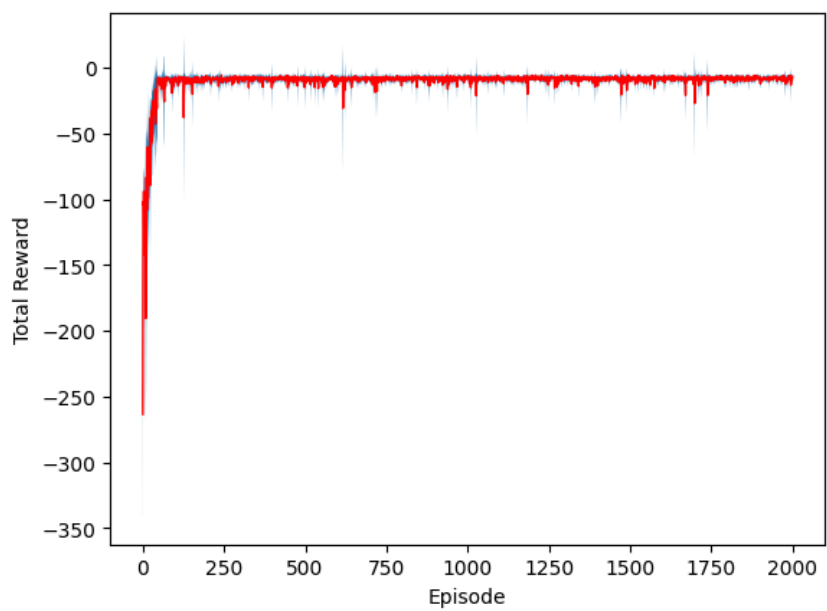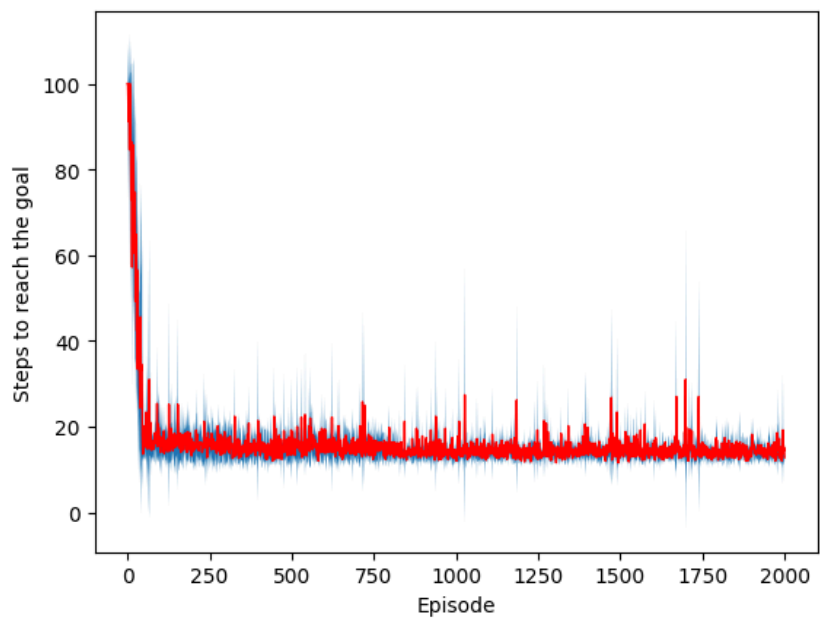Q plot



Average State Visits per episode

Because of wind, the agent converges to the goal state closer to the starting point. It does seem like it explores the environment a lot.
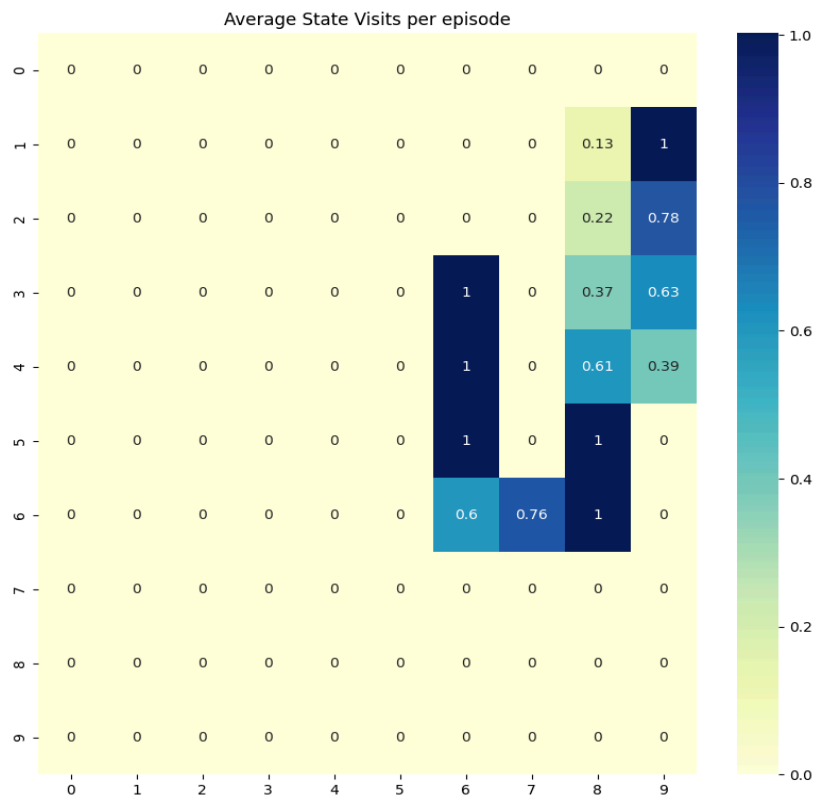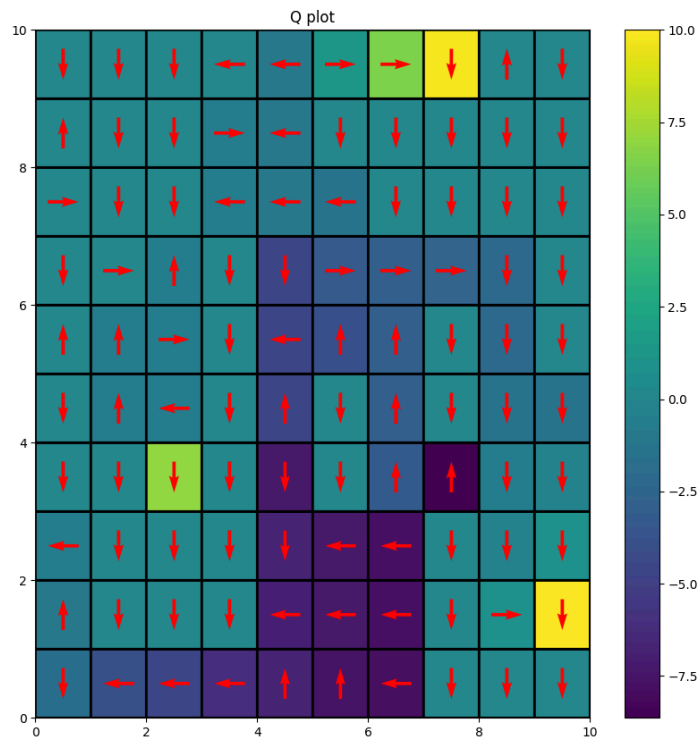
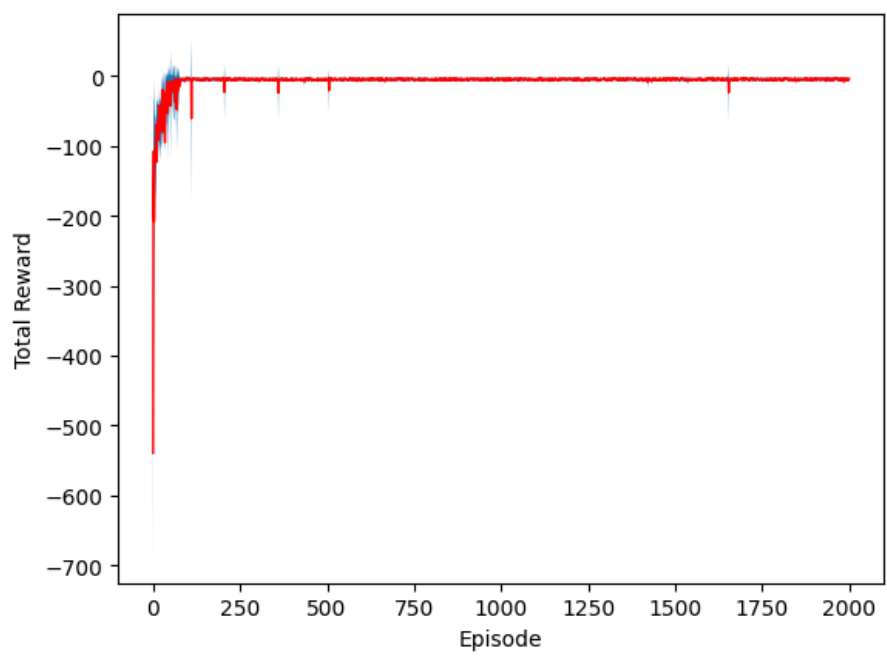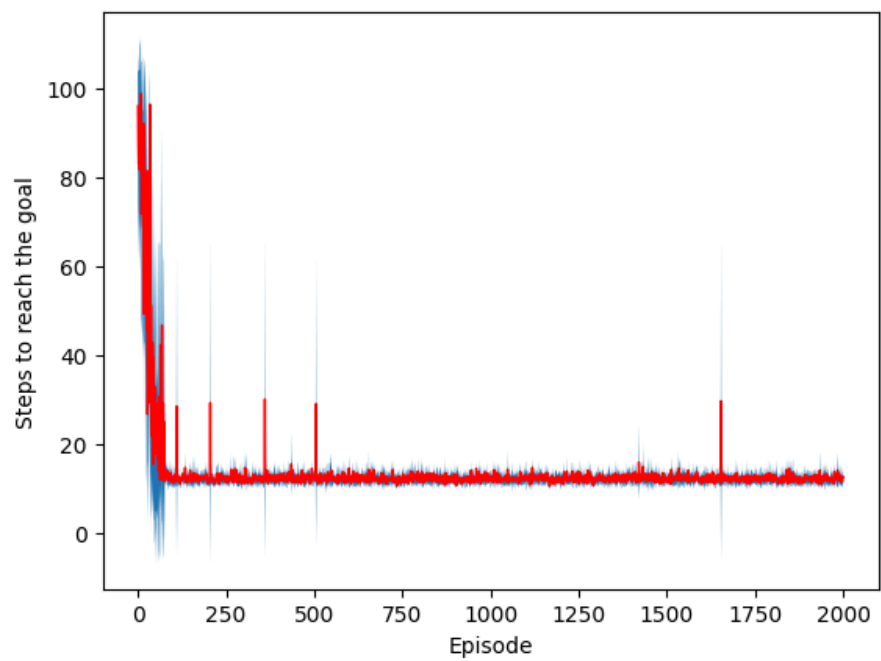Algorithm converges to around mean total reward -5 and steps at 17.

**10. Start state => [0,4] ; Wind => True; P_transition = 1 ; Policy = E-greedy:**


Q plot
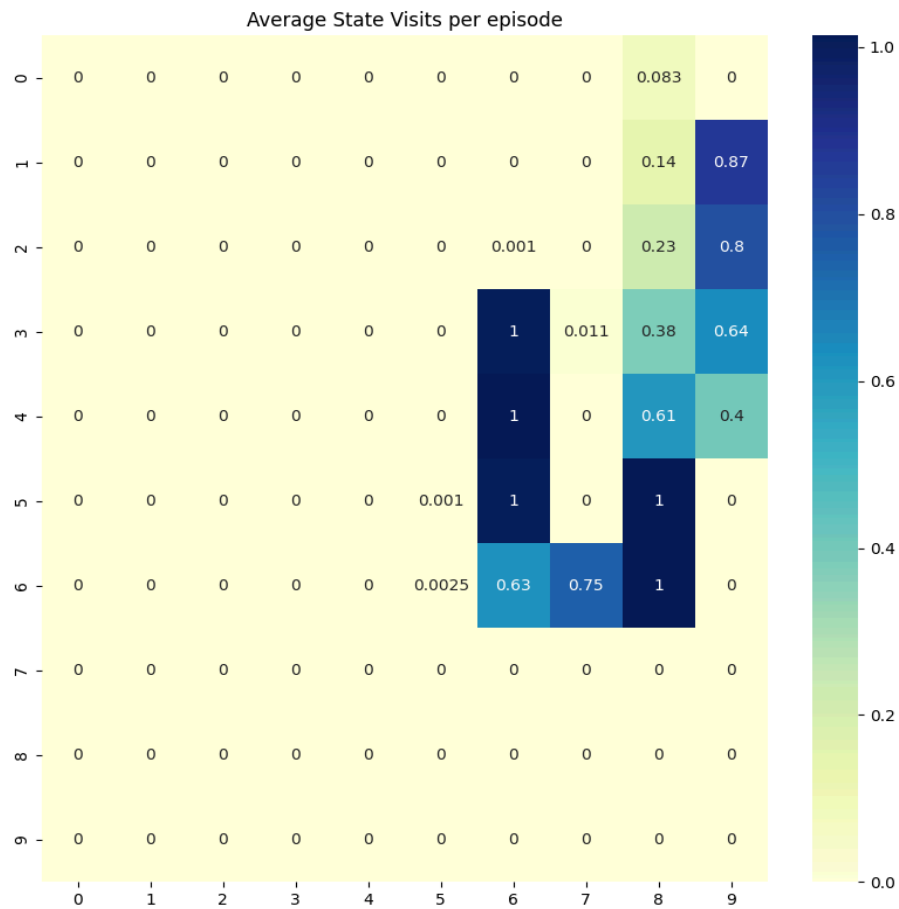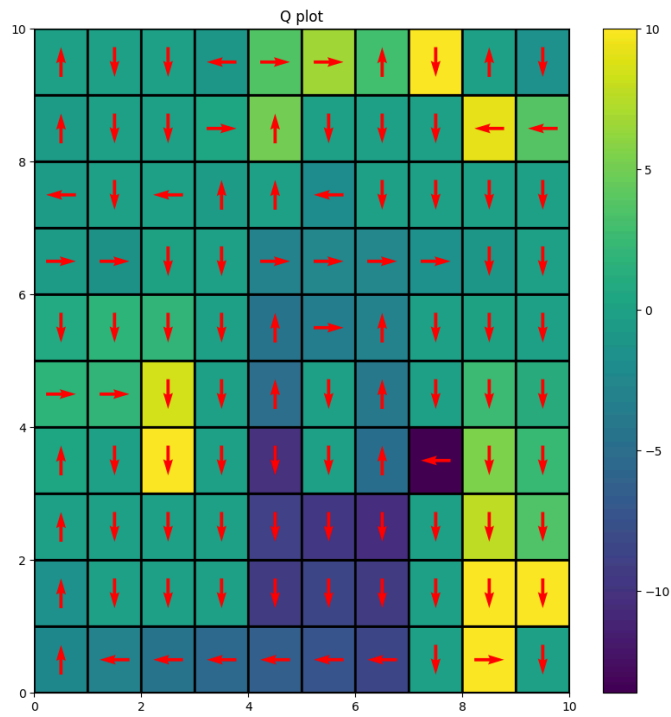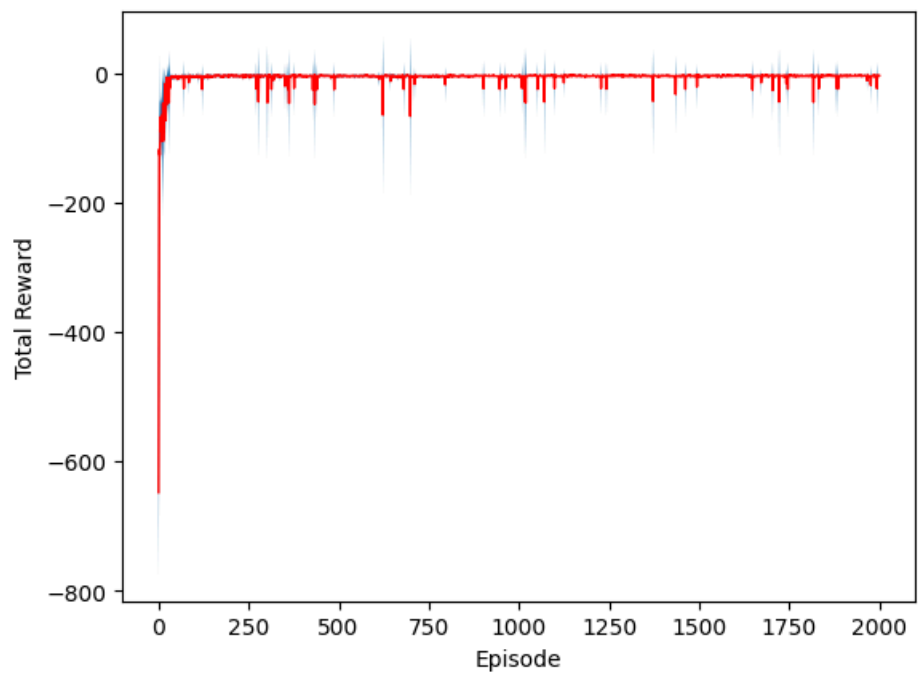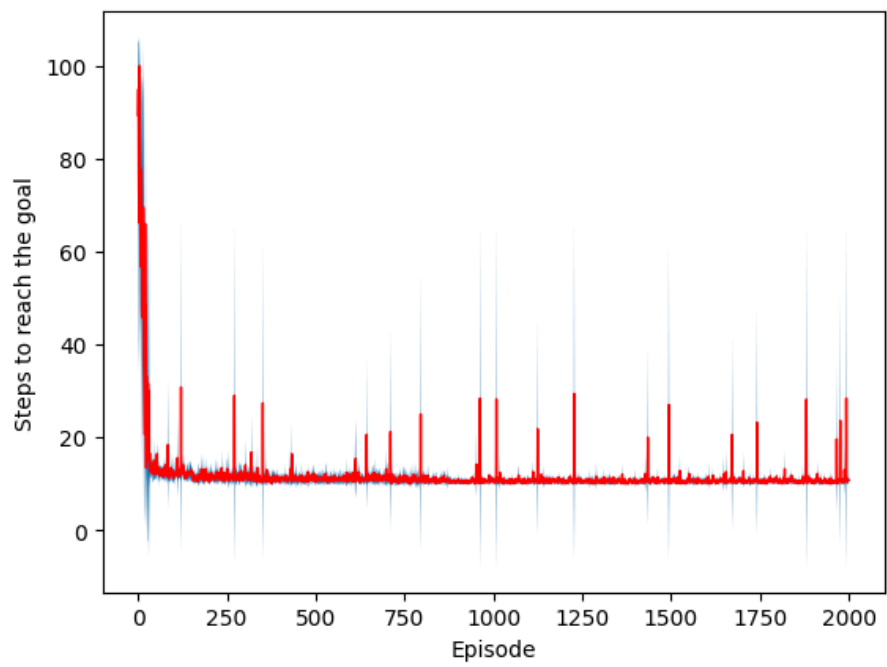

Average State Visits per episode

## 11. Start state => [3,6] ; Wind => True ; P_transition = 1 ; Policy = Softmax :



Q plot



Average State Visits per episode

**12. Start state => [3,6] ; Wind => True ; P_transition = 1 ; Policy = E-greedy:**



Q plot



Average State Visits per episode

**Comparison SARSA vs Q-learning :**

1. It seems like both SARSA and Q-learning learn similar or same paths to the same goal state.
2. But it seems like, in some cases, Q-learning is able to find the best possible path to the best goal state.
3. However in case of stochasticity in the environment, SARSA opts to find a safer path that avoids restart states.