# INDIAN INSTITUTE OF TECHNOLOGY, MADRAS

## REINFORCEMENT LEARNING PROGRAMMING ASSIGNMENT 2

CS6700

---

# SARSA,Qlearning,Policy Gradient

---

*Author*
DODIYA KEVAL

*Roll Number*
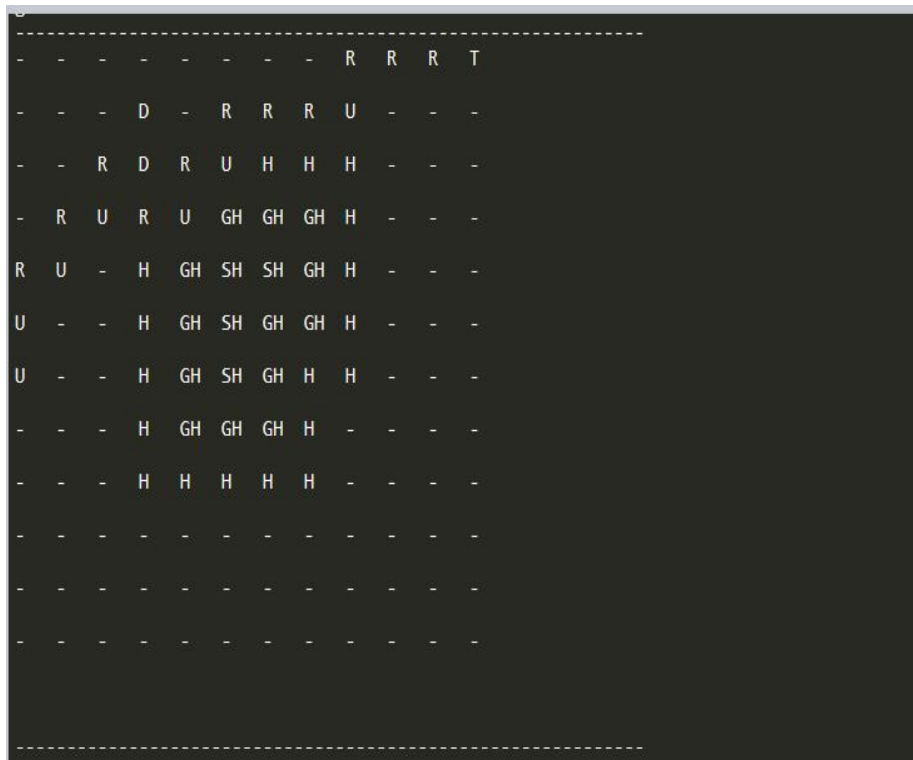CS19M023

March 15, 2020

# 1 Problem 1

## 1.1 Q Learning



Figure 1: Qlearning optimal policy for goal A

here i'm assuming H = hole with reward -1,GH = hole with reward -2, SH = hole with reward -3, and U = upper direction, L = left, R = right, D = down,T = terminal and '-' = empty slot.
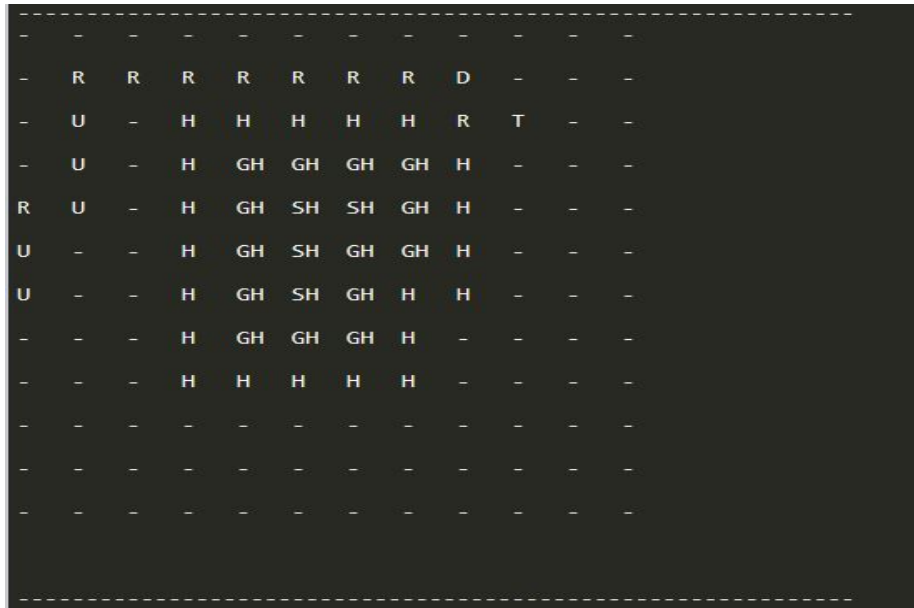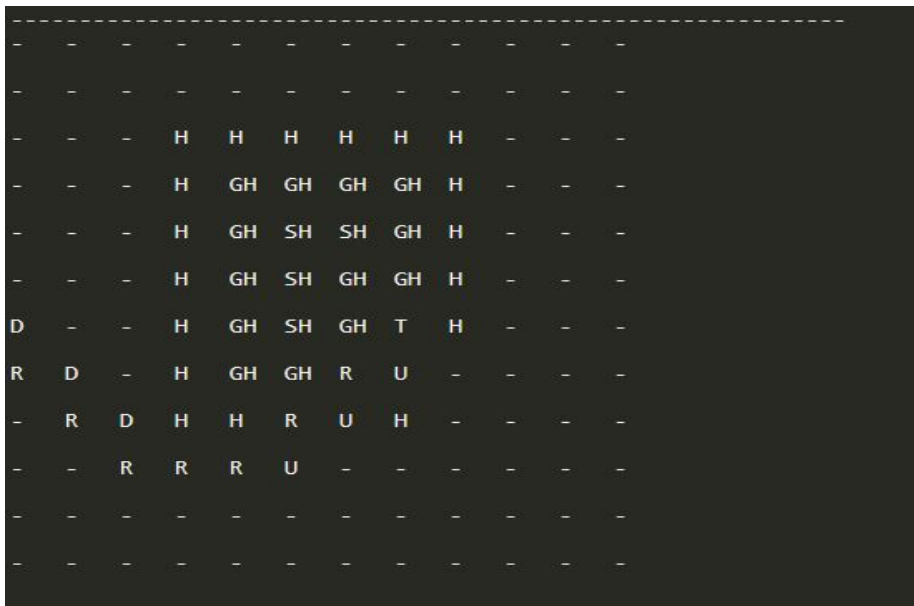
Figure 2: Qlearning optimal policy for goal B
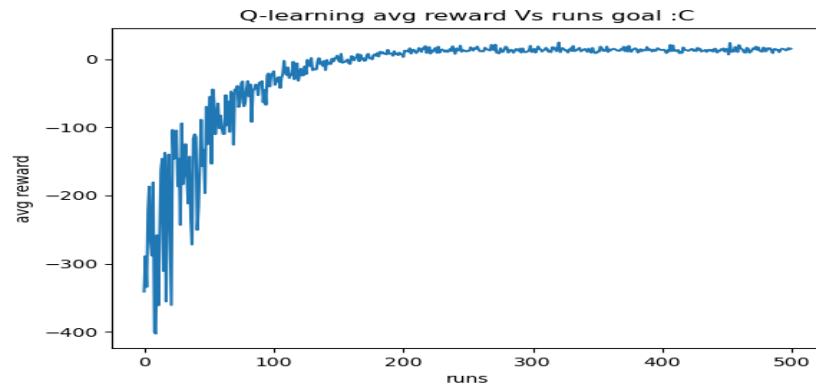


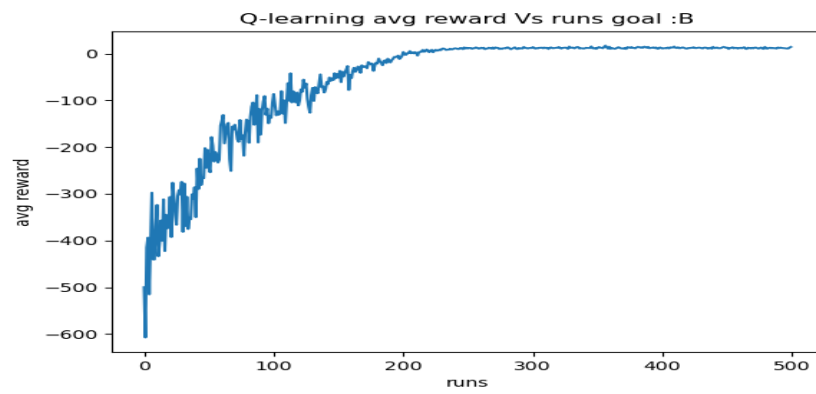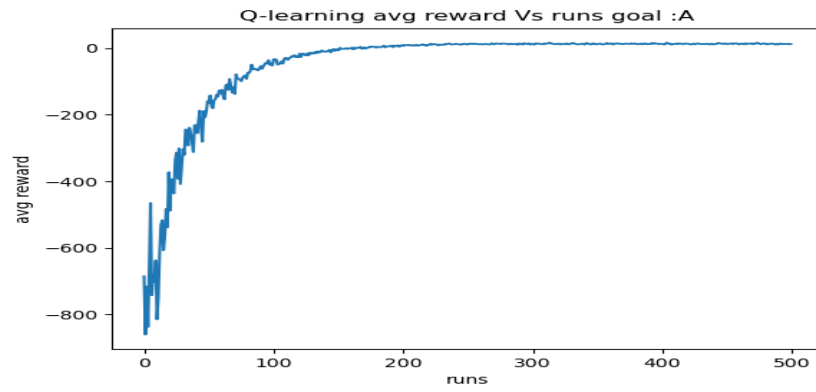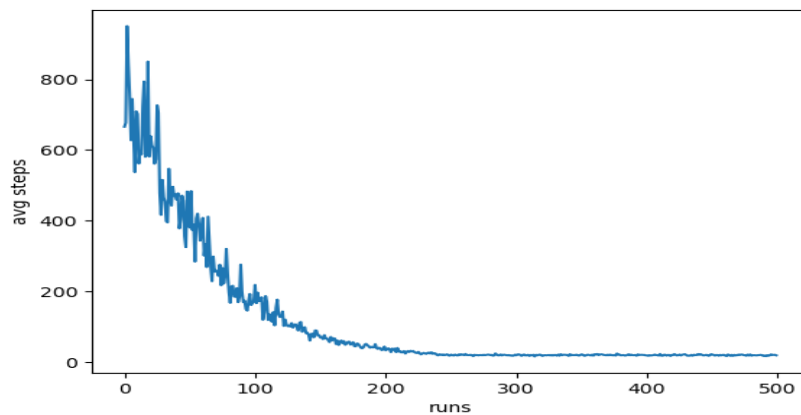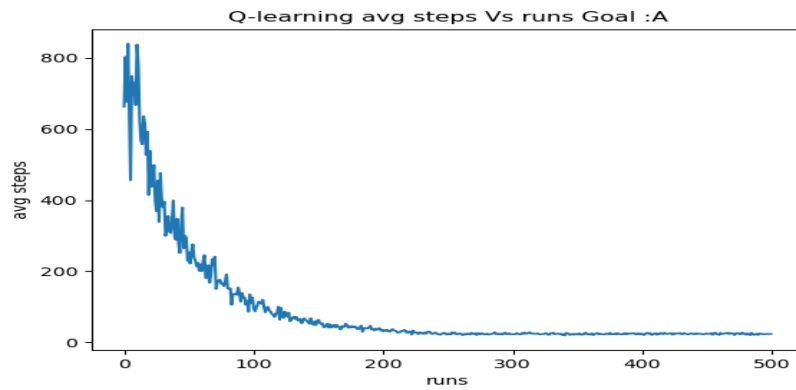Figure 3: Qlearning optimal policy for goal C

Figure 4: QLearning Reward VS runs with different Goals

- I counted offgrid Reward as -5 and 0 reward for normal action hence minimum possible reward goes very low. and i implemented a e-greedy with 1 as initial epsilon so that algorithm explore more in intial stage then it slowly decreases to 0. so i'm increasing exploitation steps.

- for Q learning initially it explores and eventually it learns optimal policy hence graph of reward increases as shows in graph.

- for initial steps variance in reward is more because of exploration and then after it converges(finding best policy of states) to best rewards

- i ran 50 independent experiment with 500 episodic task and again i register only one gridworld with different goals.

- showing policy is optimal policy for all runs
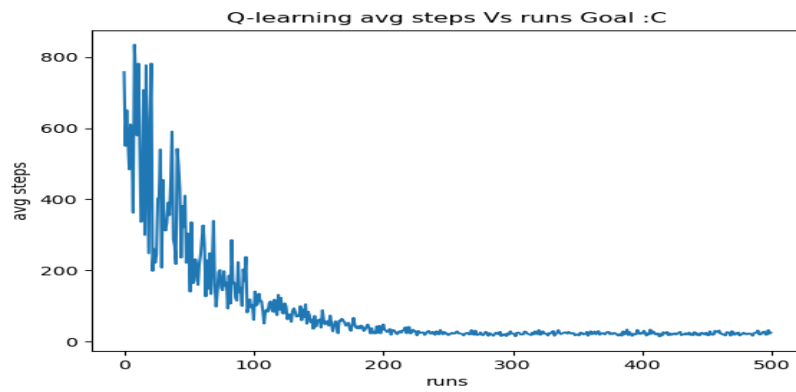


Q-learning avg steps Vs runs Goal :A

Figure 5: QLearning steps VS runs with different Goals

- here initially i got maximum steps 800 around then with number of runs increases that eventually steps decreases.

- and variance of goal C is more that is because of random behaviour and exploration of Q learning is more hence its take more steps
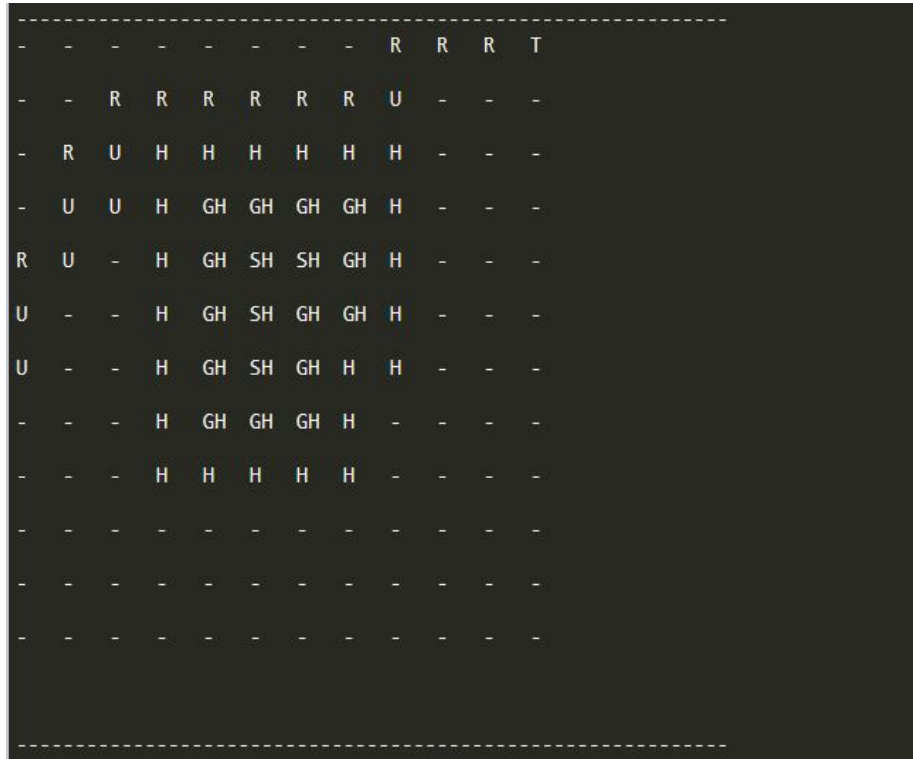
## 1.2 SARSA :
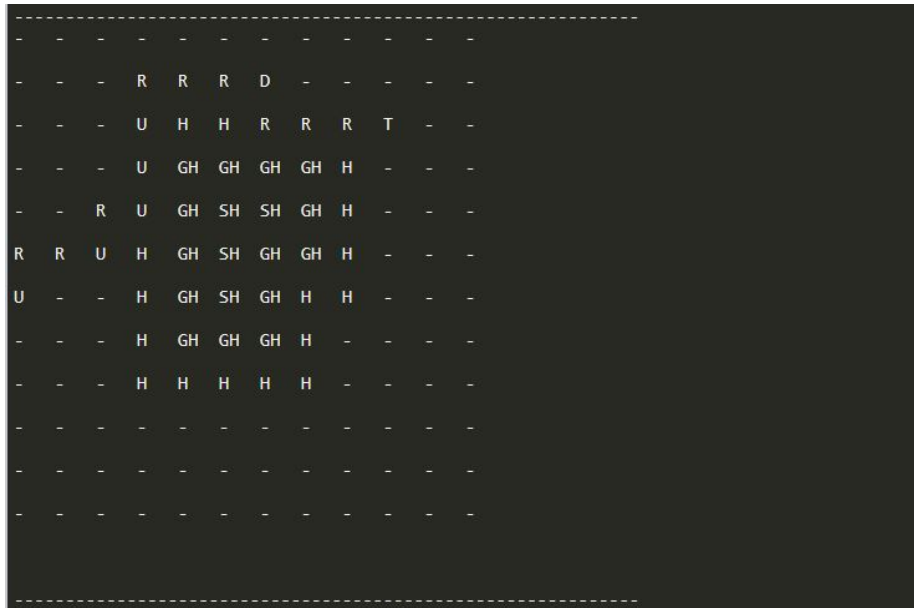


Figure 6: SARSA optimal policy for goal A
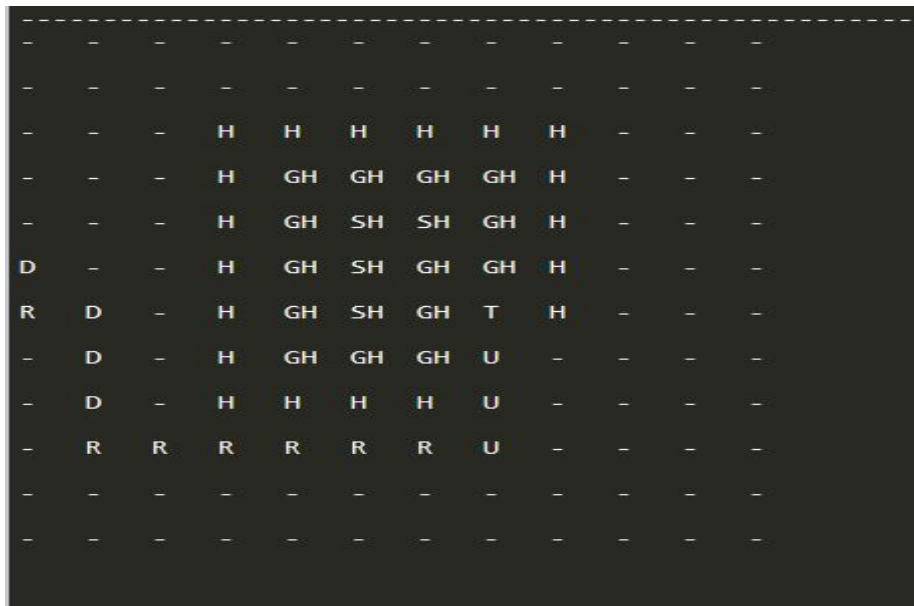
Figure 7: SARSA optimal policy for goal B



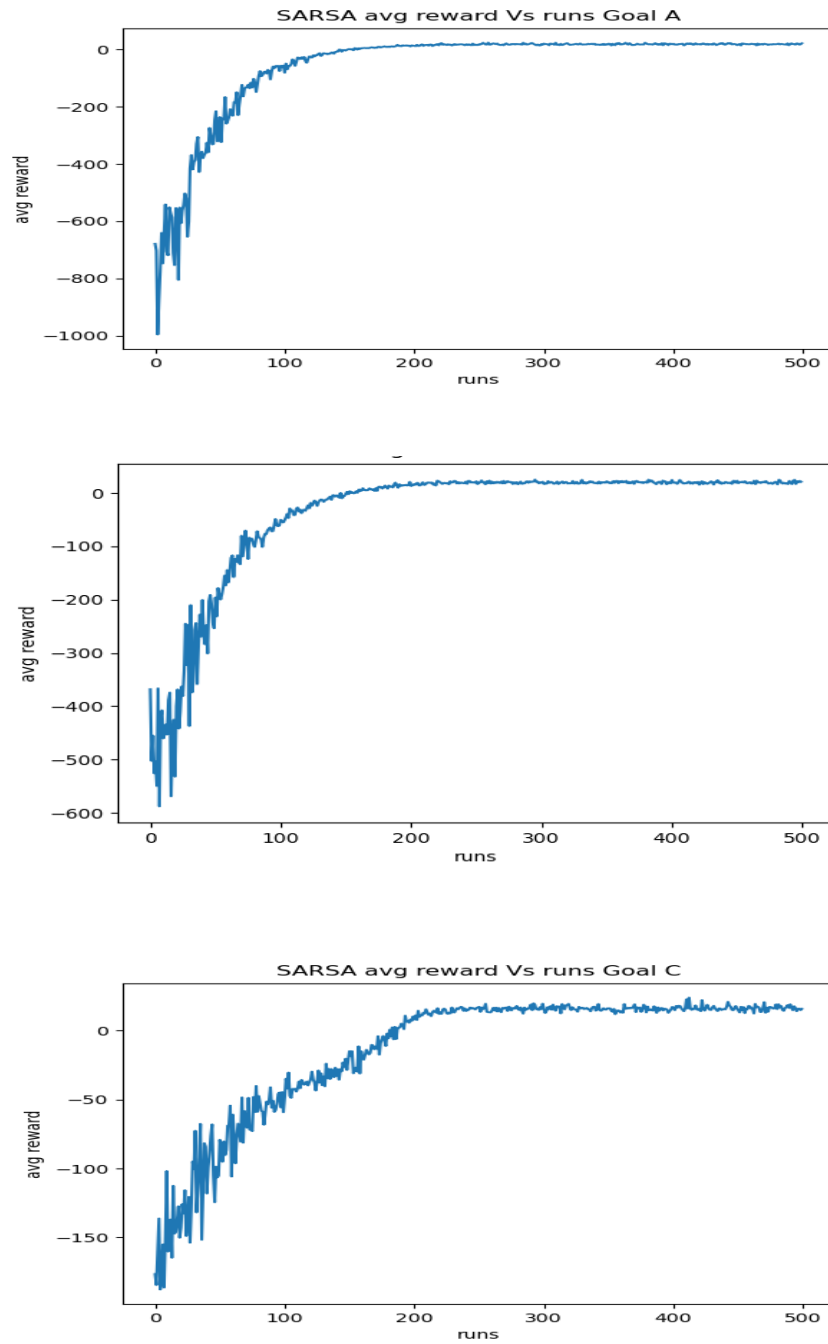Figure 8: SARSA optimal policy for goal C

Figure 9: SARSA avg Reward VS runs with different Goals with runs of 500
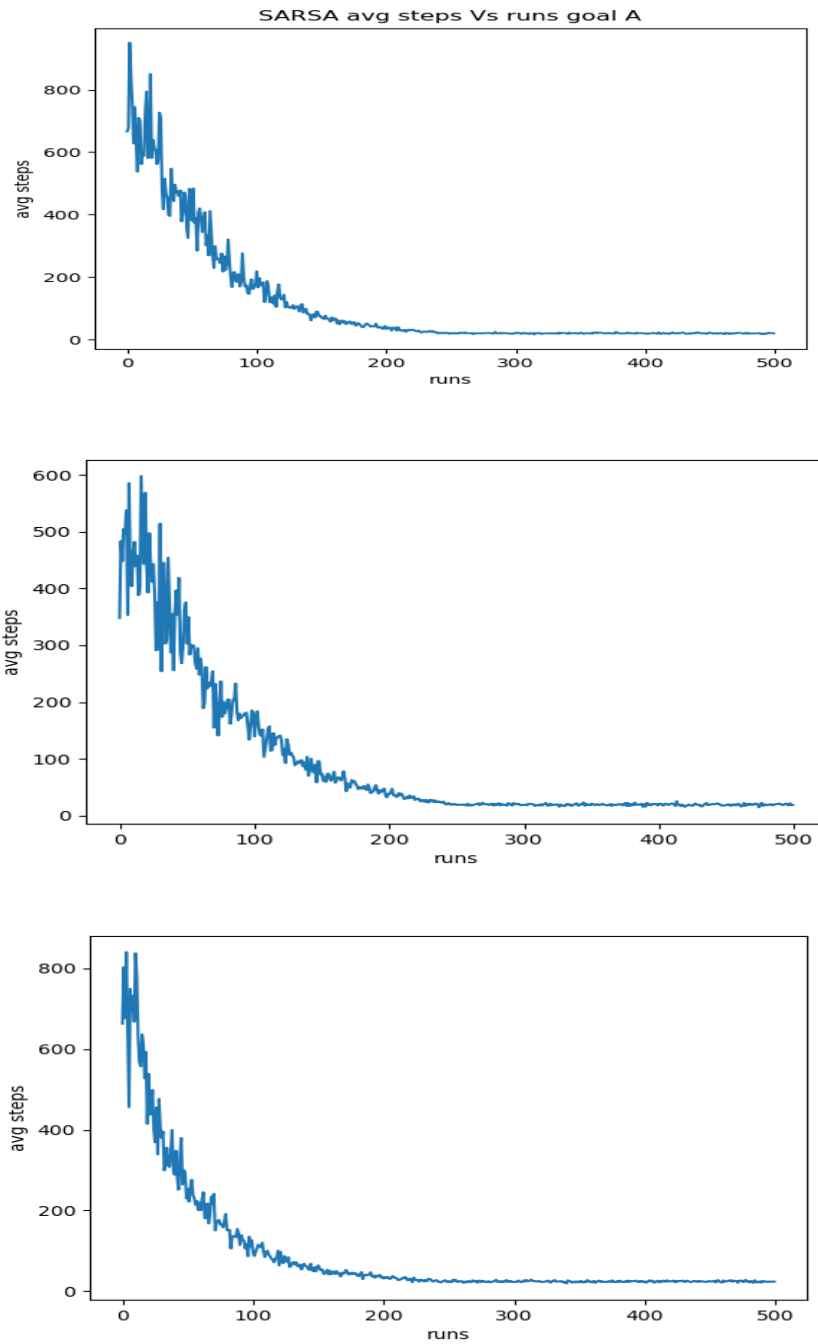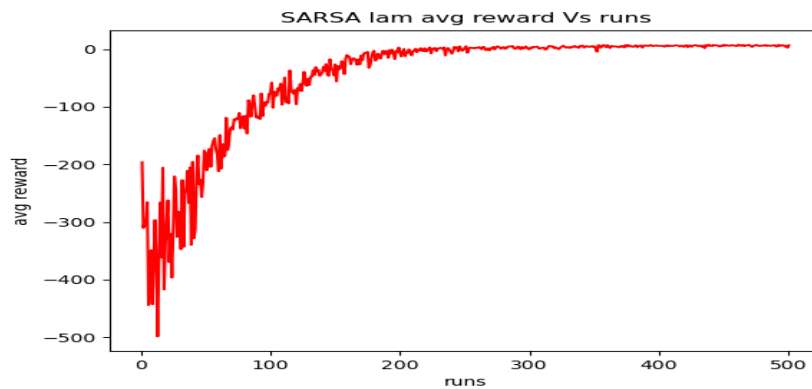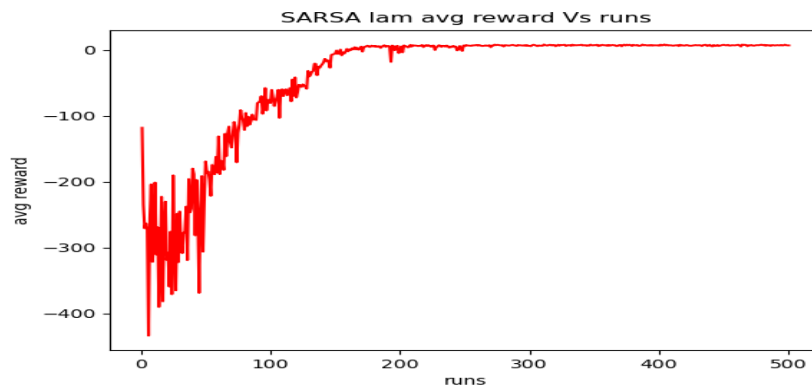and avg over 50 experiments

Figure 10: SARSA steps VS runs with different Goals first graph is for goal A, second graph is for goal B,and then goal C
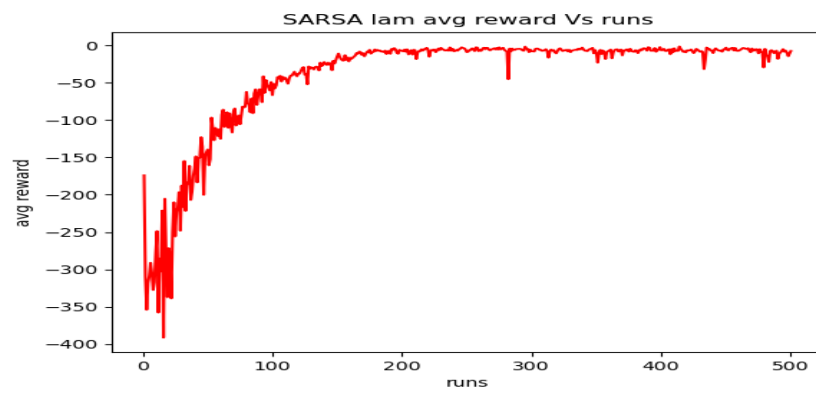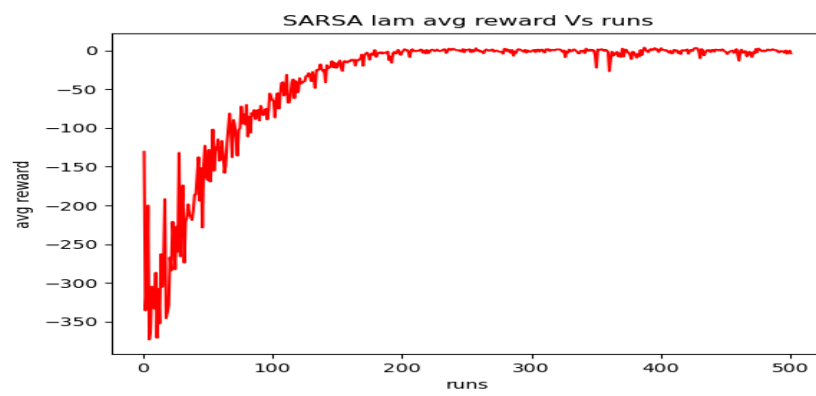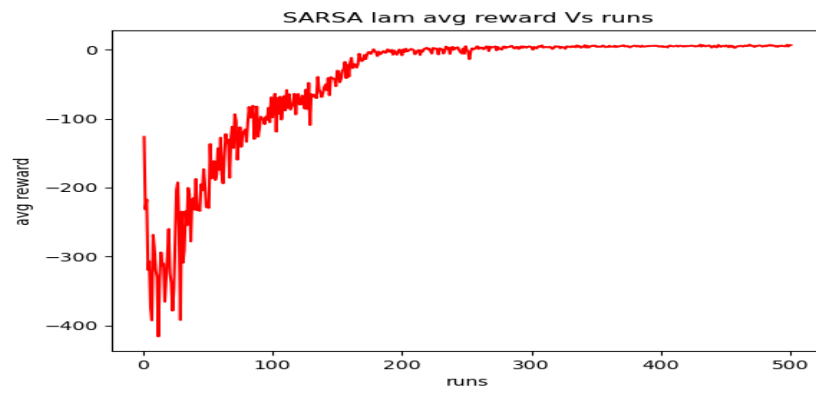
- SARSA uses its own policy to select next action's action state value hence SARSA is safer than Q learning. i put offgrid reward as -5 hence SARSA and Q learning both performing like a same algorithm with minor difference.

## 1.3 SARSA $\lambda$:

### 1.3.1 Goal A:

- following plots are of avg reward VS episodes(runs) of SARSA $\lambda$ with different $\lambda$s. and we can observed that rewards are near to same.

- epsilon value starting from 0.7 means more exploration then i'm decreasing epsilon(increment of exploitation ) with increment of episodes. hence reward goes tens to zero and above not in range 8-10.
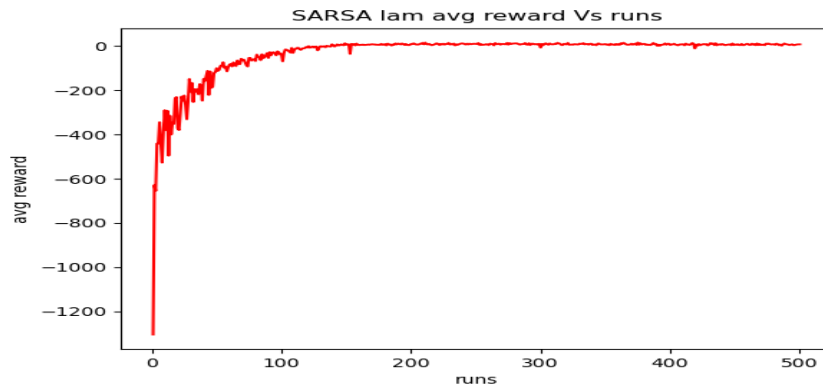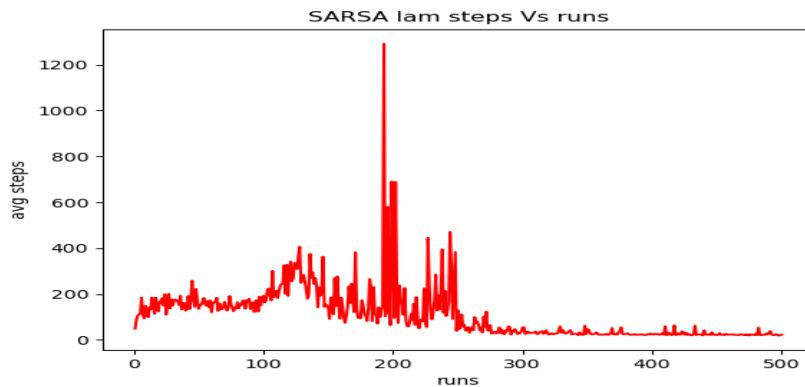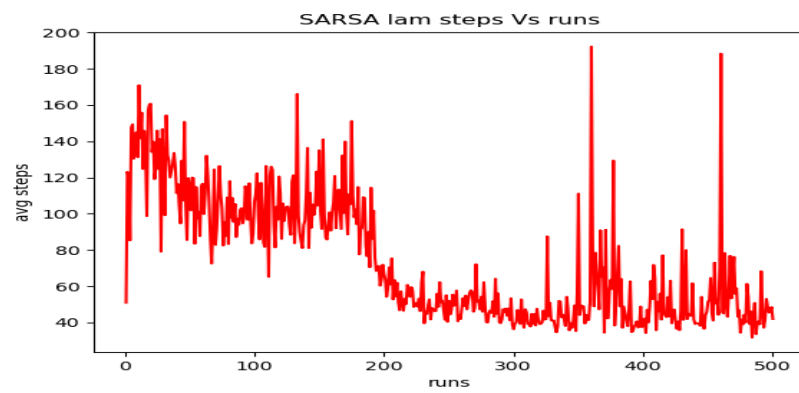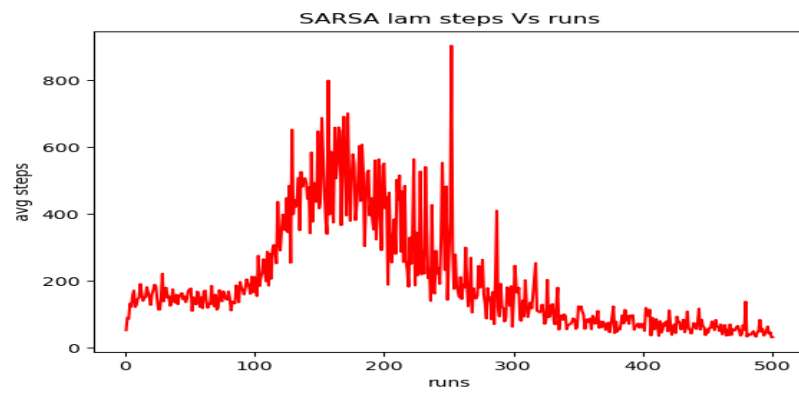
SARSA lam avg reward Vs runs



SARSA lam avg reward Vs runs
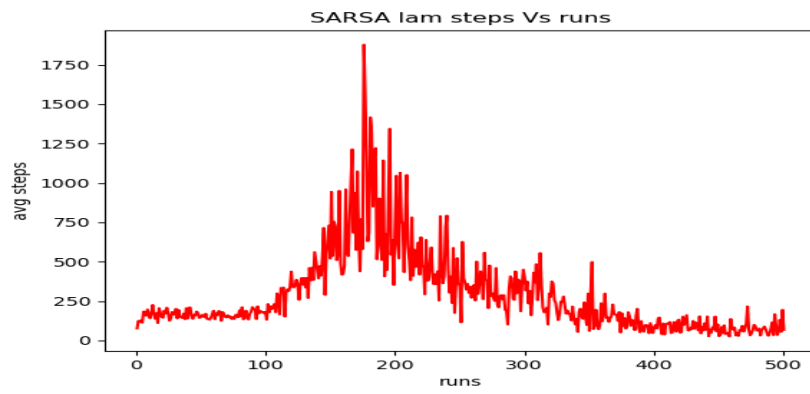


SARSA lam avg reward Vs runs

11

Figure 11: SARSA lambada reward VS episodes with different lambada over 500 episodes and 50 runs with goal A where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$

- in SARSA $\lambda$ whenever lambada value increases we can observe(see below plots) that avg steps size also got decreases.

- all reward points are converges around 500 episodes.

- I have taken offgrid reward as -5 hence minimum reward amount is very low(around -400).

- following plots are of avg steps VS episodes(runs) of SARSA $\lambda$ with different $\lambda$s. and we can observed that avg steps are decreasing with $\lambda$

- In steps Vs episodes graph we can see more variance it's because of exploration. exploration may take more steps that is varies in all runs.
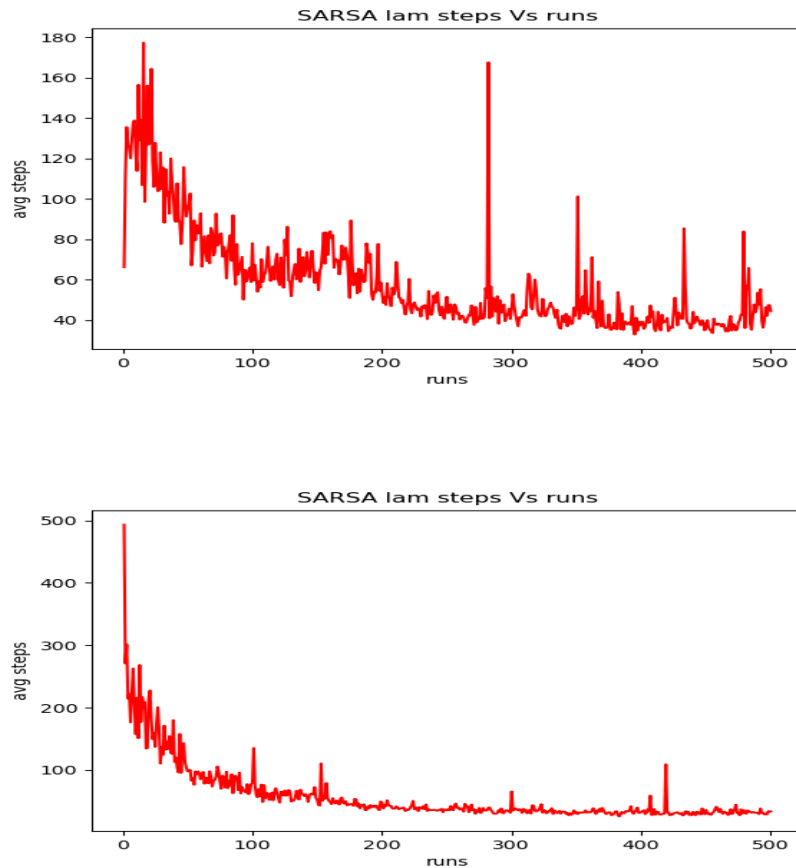
SARSA lam steps Vs runs



SARSA lam steps Vs runs
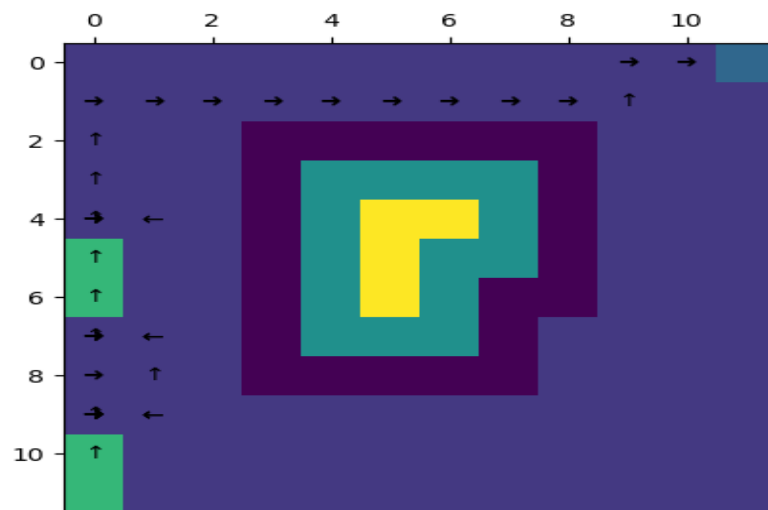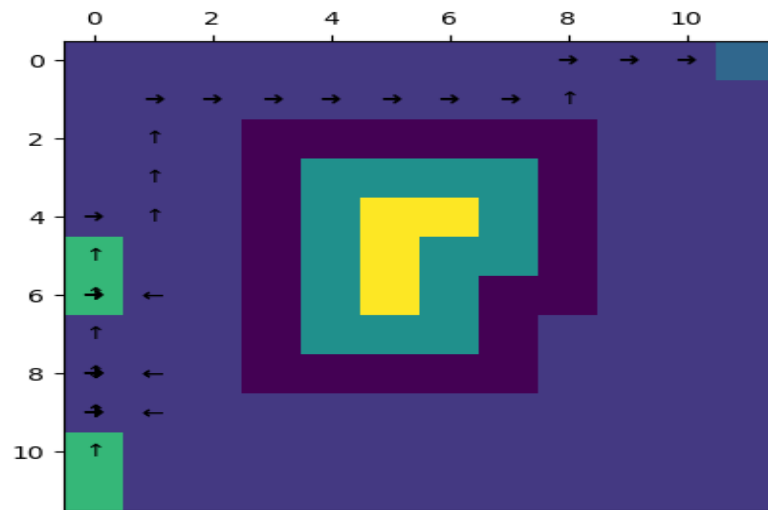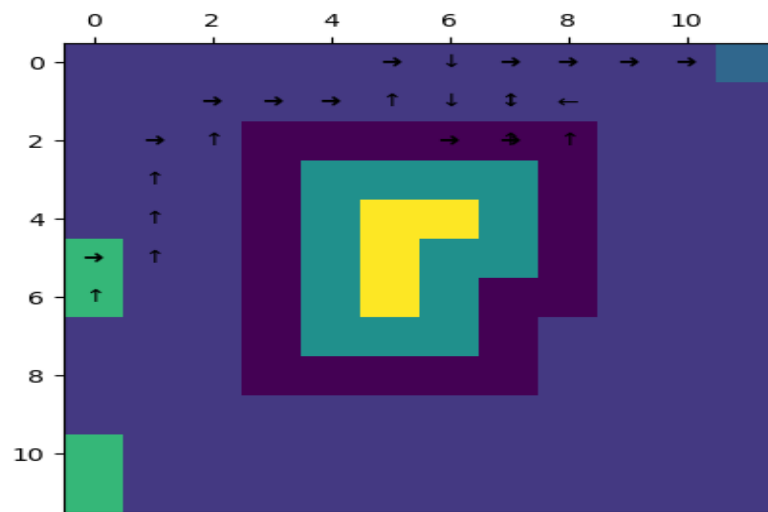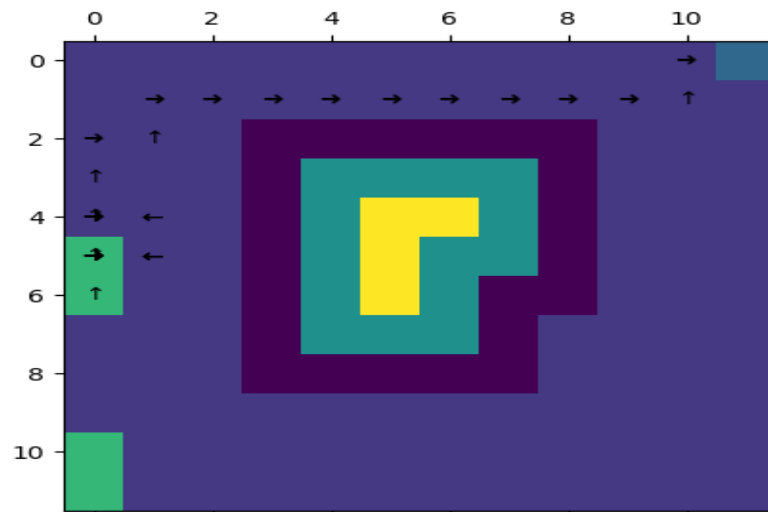


SARSA lam steps Vs runs

13

Figure 12: SARSA lambada avg steps VS episodes with different lambada over 500 episodes and 50 runs with goal A where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$

- following policies are optimal policies i got with different $\lambda$s with 500 episodes and 50 runs. i uploaded policy which is optimal of all 50 runs with 500 episodes so depends on starting state policies may take different path with maximum reward and minimum steps.

- in below mentioned policies arrow indiactes directions and gridworld with different colors have different reward value (e.x. yellow color describes -3 reward,magenta color = -1,skyblue = -2,dark blue = 10(terminal) ,light green = strating states)

- some gridbox have more than one direction discribes more than one actions in same gridbox when it return back to same state.
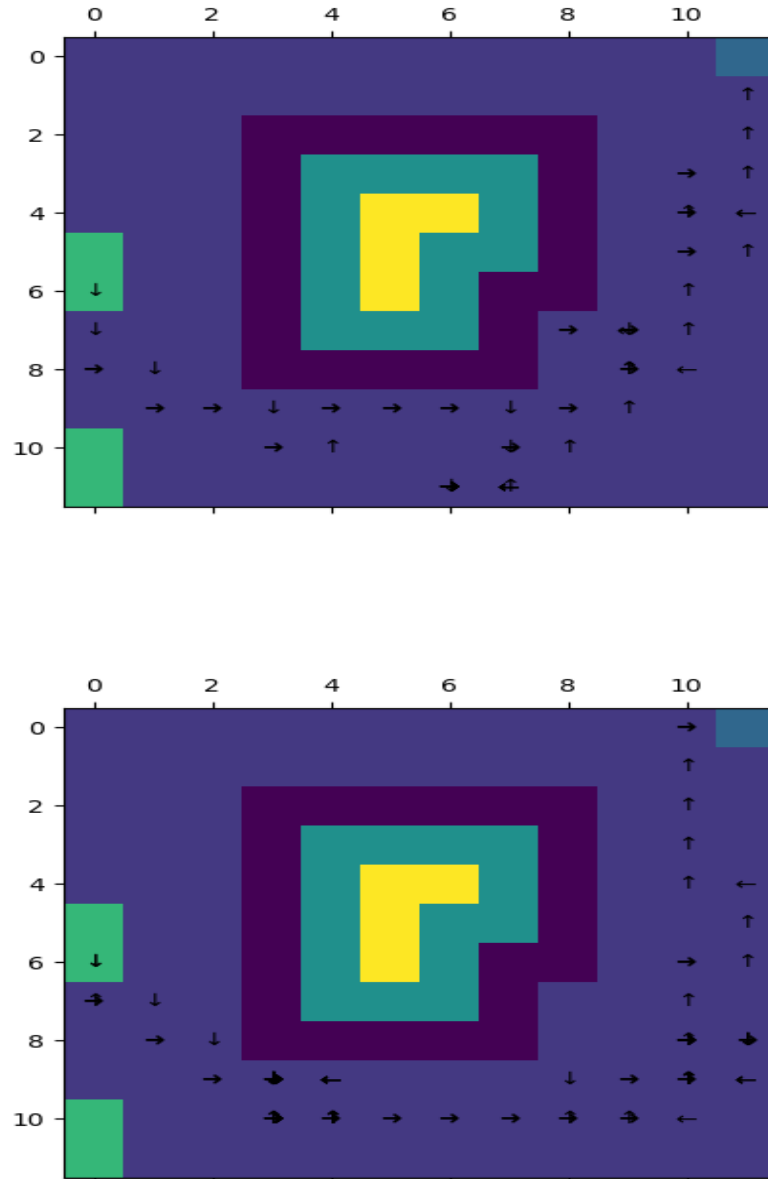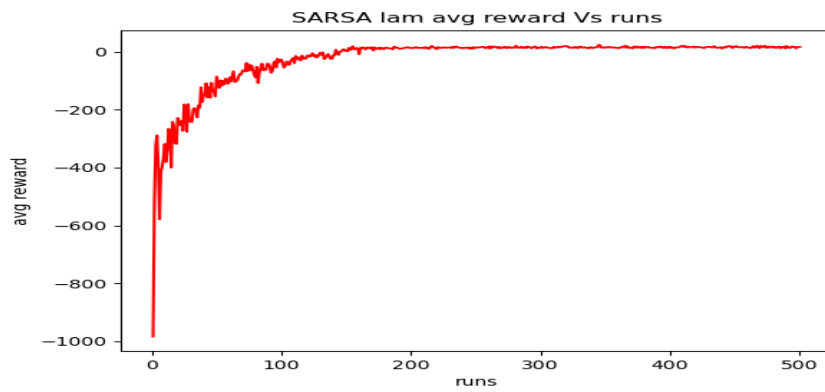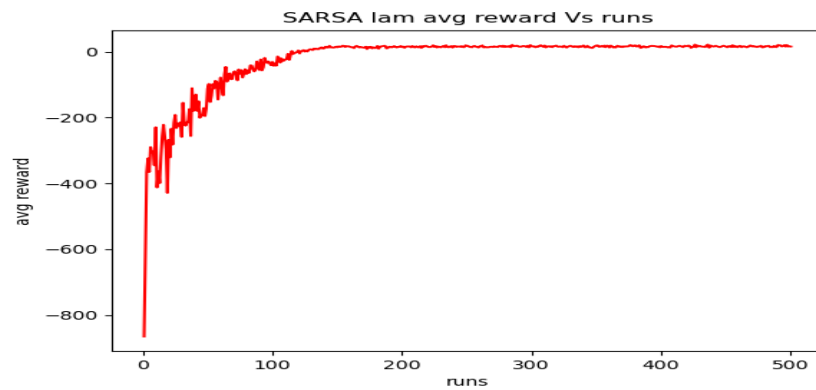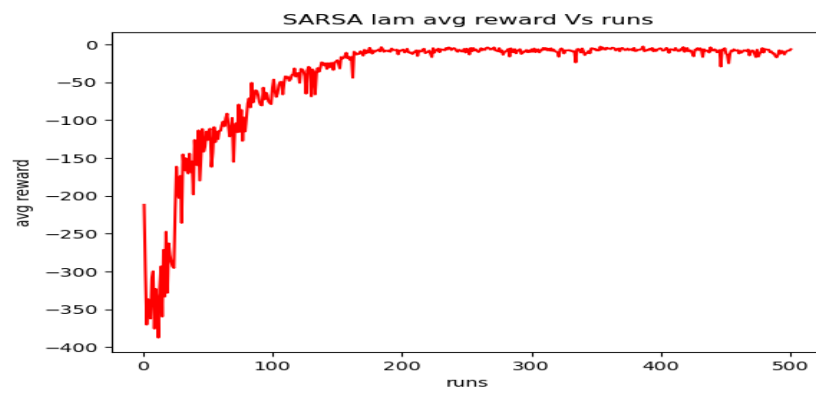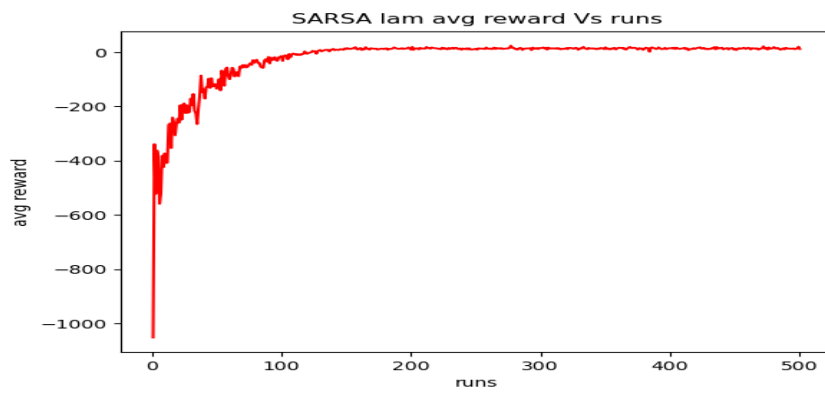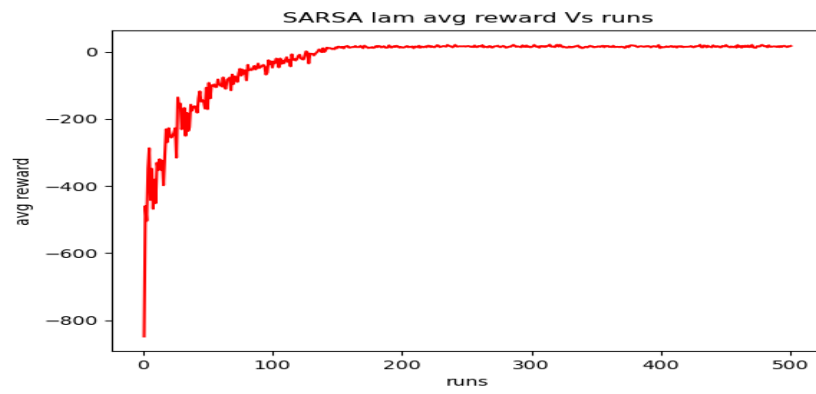
14

Figure 13: SARSA lambada optimal policies with different lambada over 500 episodes and 50 runs with goal A .where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$

### 1.3.2 Goal B:

- following plots are of avg reward VS episodes(runs) of SARSA $\lambda$ with different $\lambda$s. and we can observed that rewards are near to same.

- epsilon value starting from 0.7 means more exploration then i'm decreasing epsilon(increment of exploitation ) with increment of episodes. hence reward goes tens to zero and above not in range 8-10.
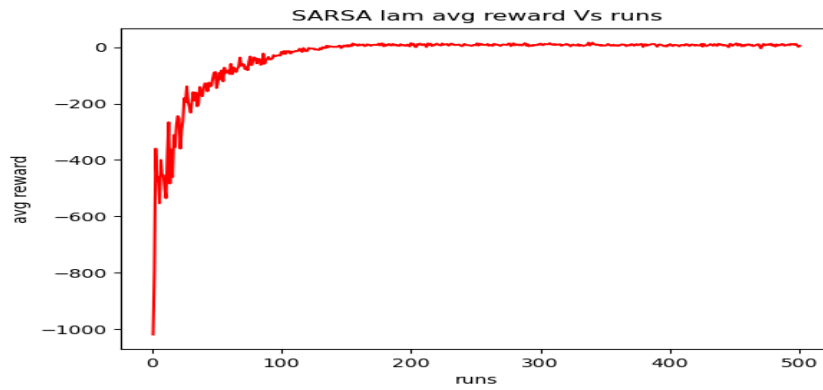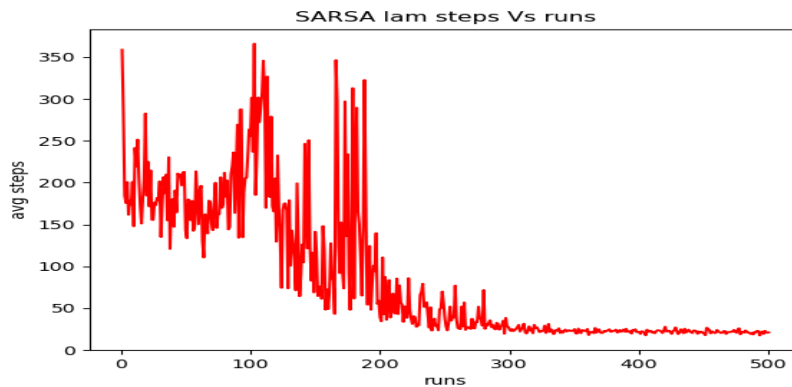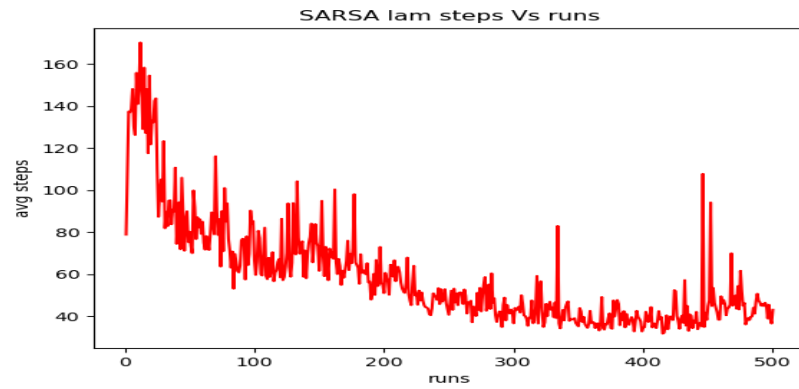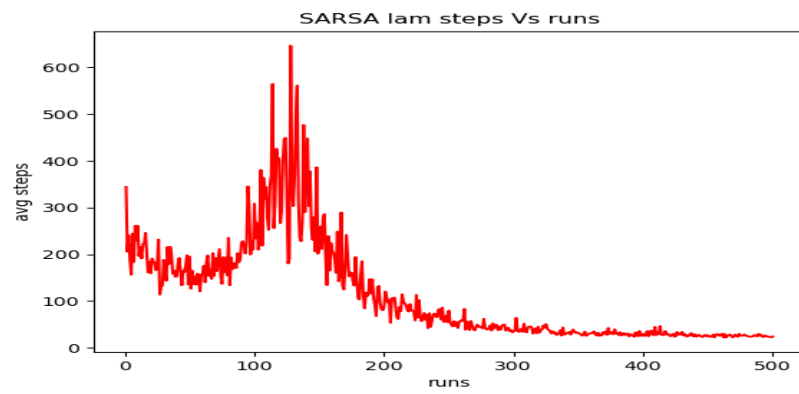
SARSA lam avg reward Vs runs



SARSA lam avg reward Vs runs
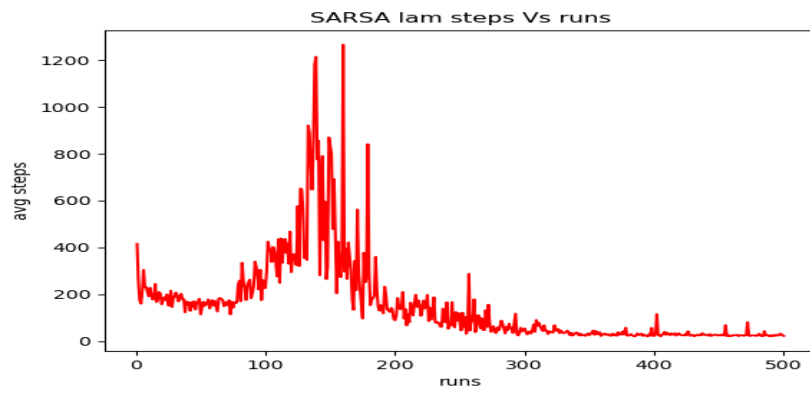


SARSA lam avg reward Vs runs

19

Figure 14: SARSA lambada reward VS episodes with different lambada over 500 episodes and 50 runs with goal B where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$

- in SARSA $\lambda$ whenever lambada value increases we can observe(see below plots) that avg steps size also got decreases.

- all reward points are converges around 500 episodes.

- I have taken offgrid reward as -5 hence minimum reward amount is very low(around -400).

- following plots are of avg steps VS episodes(runs) of SARSA $\lambda$ with different $\lambda$s. and we can observed that avg steps are decreasing with $\lambda$

- In steps Vs episodes graph we can see more variance it's because of exploration. exploration may take more steps that is varies in all runs.
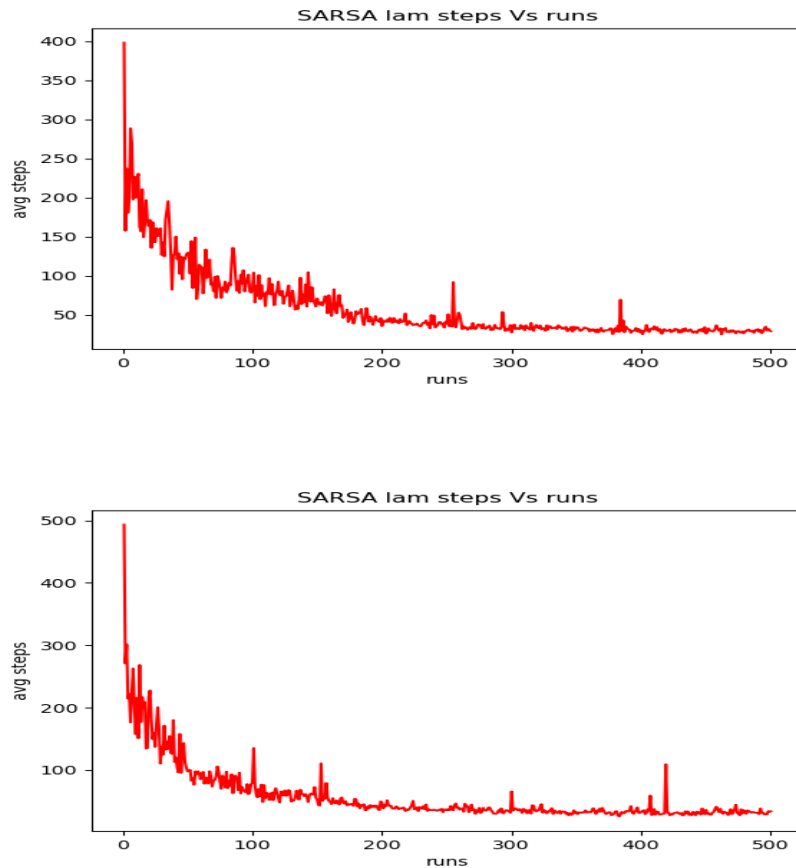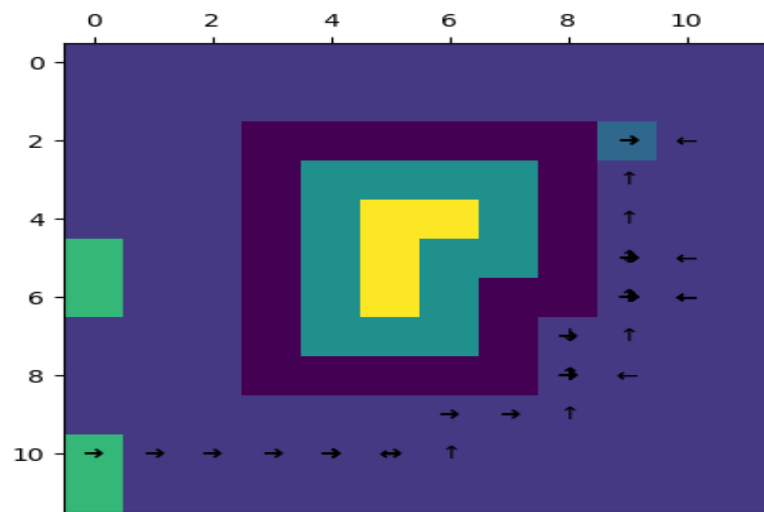


20

SARSA lam steps Vs runs

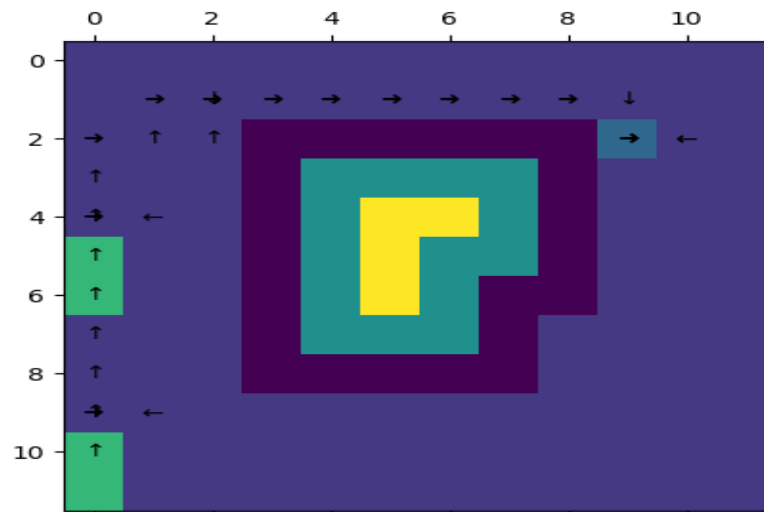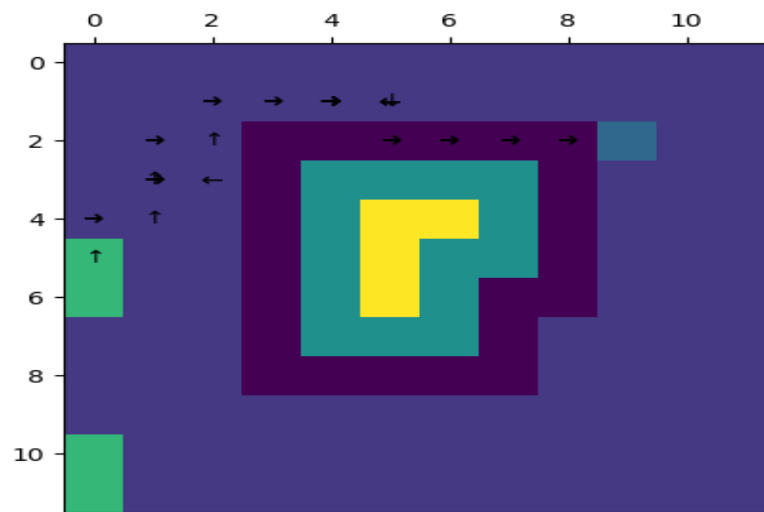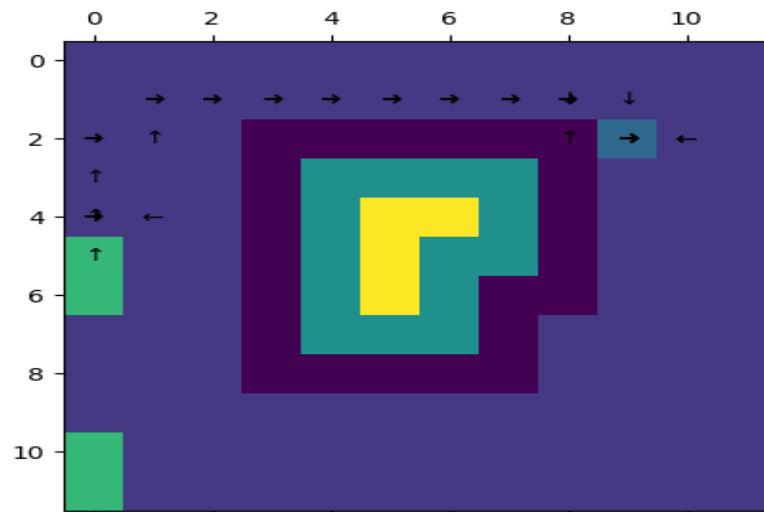SARSA lam steps Vs runs

SARSA lam steps Vs runs

Figure 15: SARSA lambada avg steps VS episodes with different lambada over 500 episodes and 50 runs with goal B where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$

- following policies are optimal policies i got with different $\lambda$s with 500 episodes and 50 runs. i uploaded policy which is optimal of all 50 runs with 500 episodes so depends on starting state policies may take different path with maximum reward and minimum steps.

- in below mentioned policies arrow indiactes directions and gridworld with different colors have different reward value (e.x. yellow color describes -3 reward,magenta color = -1,skyblue = -2,dark blue = 10(terminal) ,light green = strating states)

- some gridbox have more than one direction discribes more than one actions in same gridbox when it return back to same state.
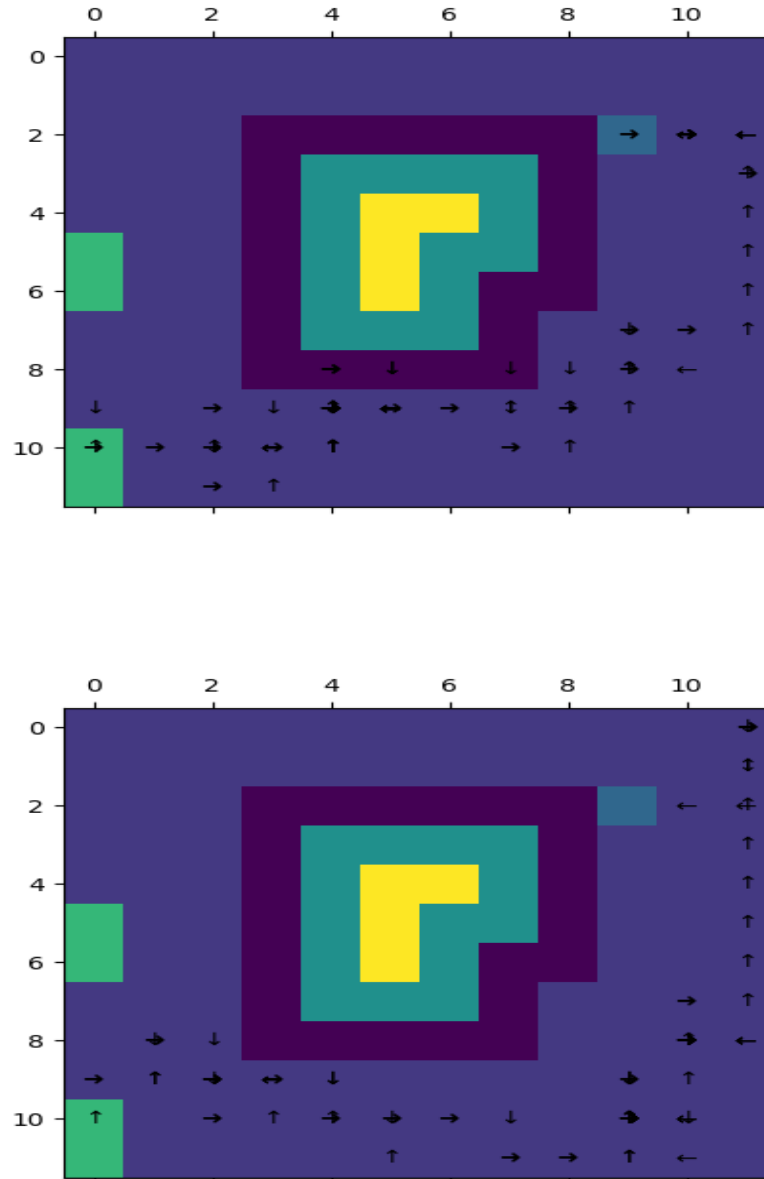
Figure 16: SARSA lambada optimal policies with different lambada over 500 episodes and 50 runs with goal B .where $\lambda = [0, 0.3, 0.5, 0.9, 0.99, 1]$
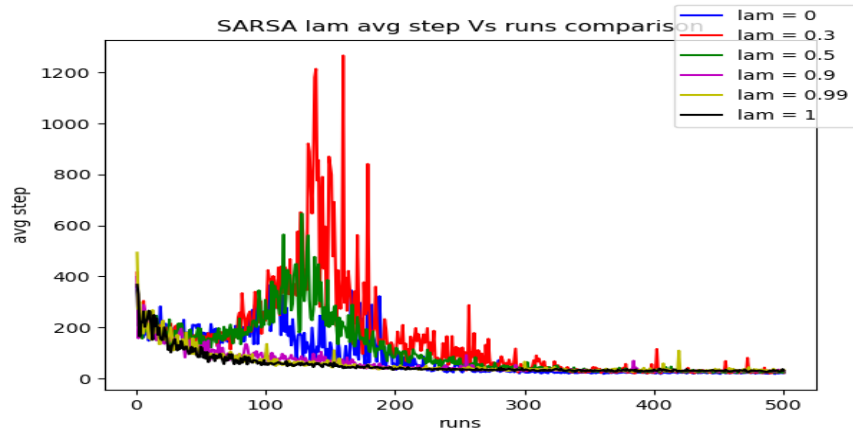
Figure 18: SARSA lambada comparison plots(avg steps VS epsiodes) with different lambada over 500 episodes and 50 runs with goal B .where $\lambda$ $=[0, 0.3, 0.5, 0.9, 0.99, 1]$
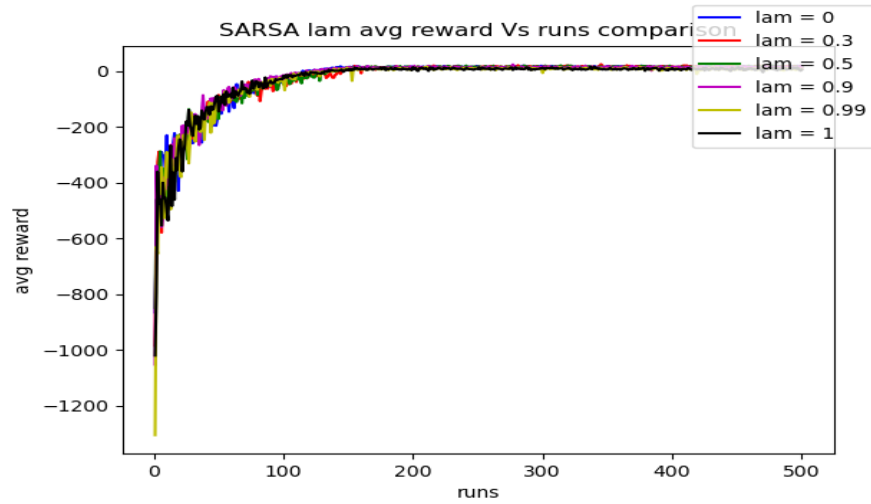


Figure 17: SARSA lambada comparison plots(avg reward VS epsiodes) with different lambada over 500 episodes and 50 runs with goal B .where $\lambda$ $=[0, 0.3, 0.5, 0.9, 0.99, 1]$