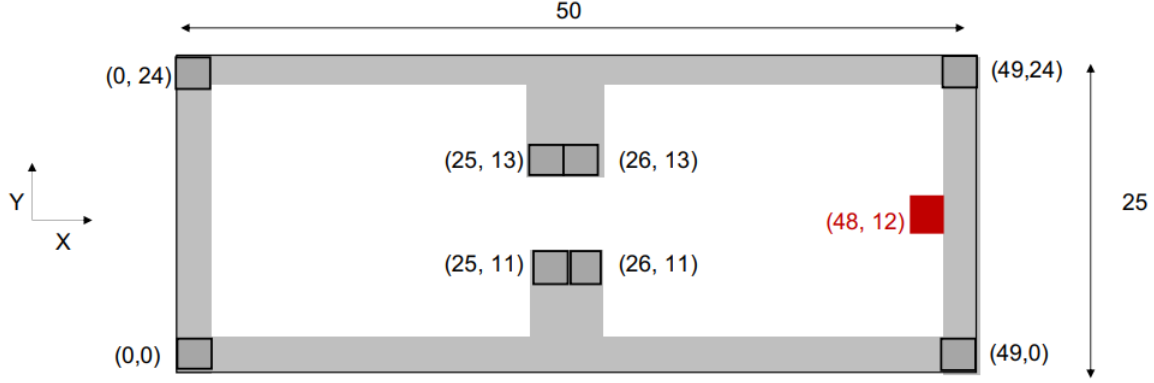


Assignment 2 Report

Students: Isha Chaudhary (2018EE30614), Ankit Kumar (2017MT10727)

1 MDP



1.1 Part a: Value Iteration

1.1.1 Background

Value Iteration is an iterative approach to solve the Bellman Equation for the State-Value function. It is primarily based on the Gauss-Seidel Method of solving simultaneous, coupled equations.

The State-Value function for a given policy, with time horizon, T, is given by:

$$V_T^\pi(s_t) = E_{s_{\tau:T}} \left[\sum_{\tau=t}^T \gamma^\tau R(s_\tau, \pi(s_\tau), s_{\tau+1}) \right]$$

Here, γ is the discount factor, and $R(s, a, s')$ is the reward function for the transition to s' from s , by taking action, a . Thus, the Bellman Equation for the State-Value Function, given a policy, is the following:

$$V(s) = E[R_{t+1} + \gamma V(s_{t+1}) | s_t = s]$$

or more elaborately:

$$V_T^\pi(s_t) = \sum_{s_{t+1}} p(s_{t+1} | s_t, \pi(s_t)) (R(s_t, \pi(s_t), s_{t+1}) + \gamma V_{T-1}^\pi(s_{t+1}))$$

The Value Iteration Algorithm finds the optimal policy by iteratively finding the State-Value function for each of the states and constituting a deterministic policy of actions which maximize the state-value function of a given state, computed from the State-value function of the state reached after action a .

The iterative procedure updates the state-value function as:

$$V(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

The state value function for each of the states is updated till the maximum difference between consecutive state-value functions falls below a pre-decided threshold.

After the state-value function converges, the deterministic policy returned, for a given state, s is:

$$\pi(s) = \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

1.1.2 Model

Action Model: Let the intended action be a

$$p(a) = 0.8$$

$$p(A - a) = 0.2/3, \text{ where } A \text{ is the action space.}$$

Reward Model:

$$R(\text{Collide with a wall}) = -1$$

$$R(\text{reach goal state}) = +100$$

$$R(\text{otherwise}) = 0$$

1.1.3 Implementation Parameters

- The state-values of the grid cells was randomly assigned some value between 0 and 1.
- The state-values of the wall cells was initialized as 0.
- Discount Factor, $\gamma = 0.1$
- Threshold, $\theta = 0.1$
- Maximum number of iterations = 100

1.1.4 Observations

The State-Value function for each of the feasible grid cells is shown in [1](#).

Here, the blue region indicates the walls in the grid world. All other grid cells can be visited by the agent receiving 0 reward. The grid cell, highlighted with red boundary is the goal state. The State-Value shown is after 100 iterations of Value Iteration.

The Policy obtained at the end of Value Iteration procedure is shown in [2](#)

1.1.5 Inference

1. The plot for the State-Value function is showing high value for the goal state and other states around it.
2. The high value of the goal state can be attributed to the fact that the goal state is a non-terminal state and a very high reward, +100, is received on staying in the goal state.
3. The corresponding policy map obtained indicates the presence of credit-assignment problem in the given instance. As the reward obtained, for a transition neither into a wall nor into the goal state, is 0, and the discount factor is very low, the value function of the states away from the goal couldn't be updated very much.
4. Low Discount factor value causes less influence of the reward at the goal state, on the states far away from the goal.

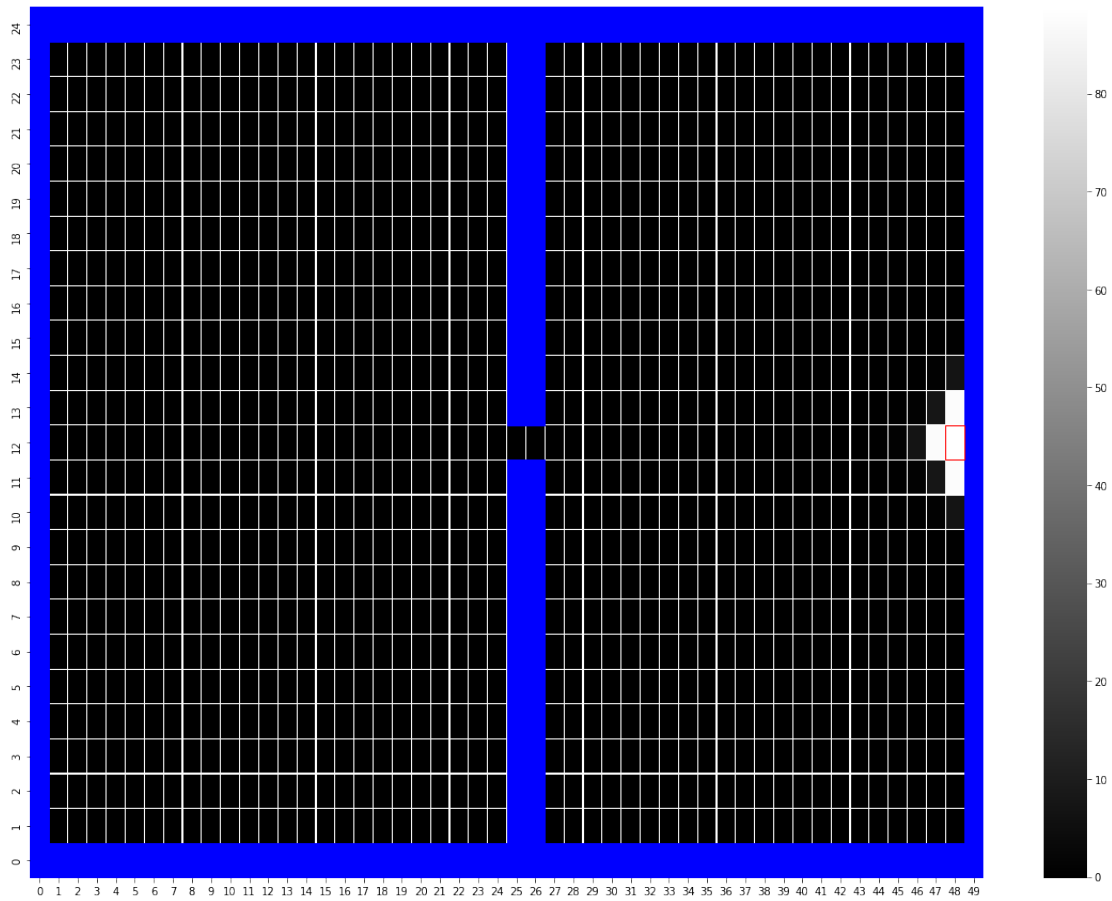


Figure 1: State Value Function

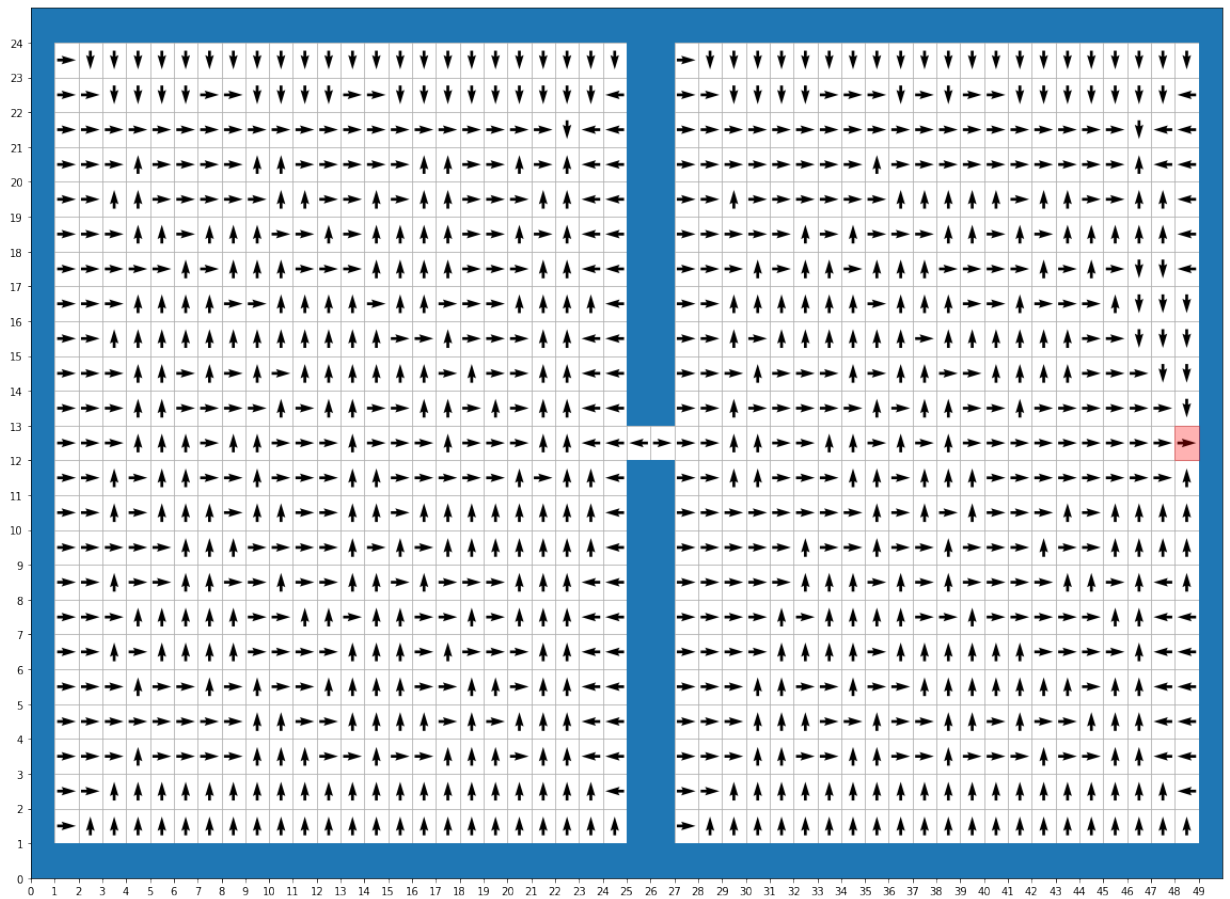


Figure 2: Optimum Policy Map obtained from Value Iteration

1.2 Part b: Varying the Discount factor

The Discount factor is changed to $\gamma = 0.99$.

1.2.1 Observations

Value Function after 20 iterations is given by the image [3](#)

Value Function after 50 iterations is given by the image [4](#)

Value Function after 100 iterations is given by the image [5](#)

In the plots, the blue boundary indicates the walls and the red outlined block is the goal state.

1.2.2 Inference

1. The discount factor used here is very high ($= 0.99$). This causes the goal state (future state) reward to contribute more to the state-value function of the grid cells away from the goal.
2. As the number of iterations of the algorithm increase, the value function converges more and more and better incorporates the effect of the future states (including the goal state). Hence the state-value of all the grid cells is seen to increase.

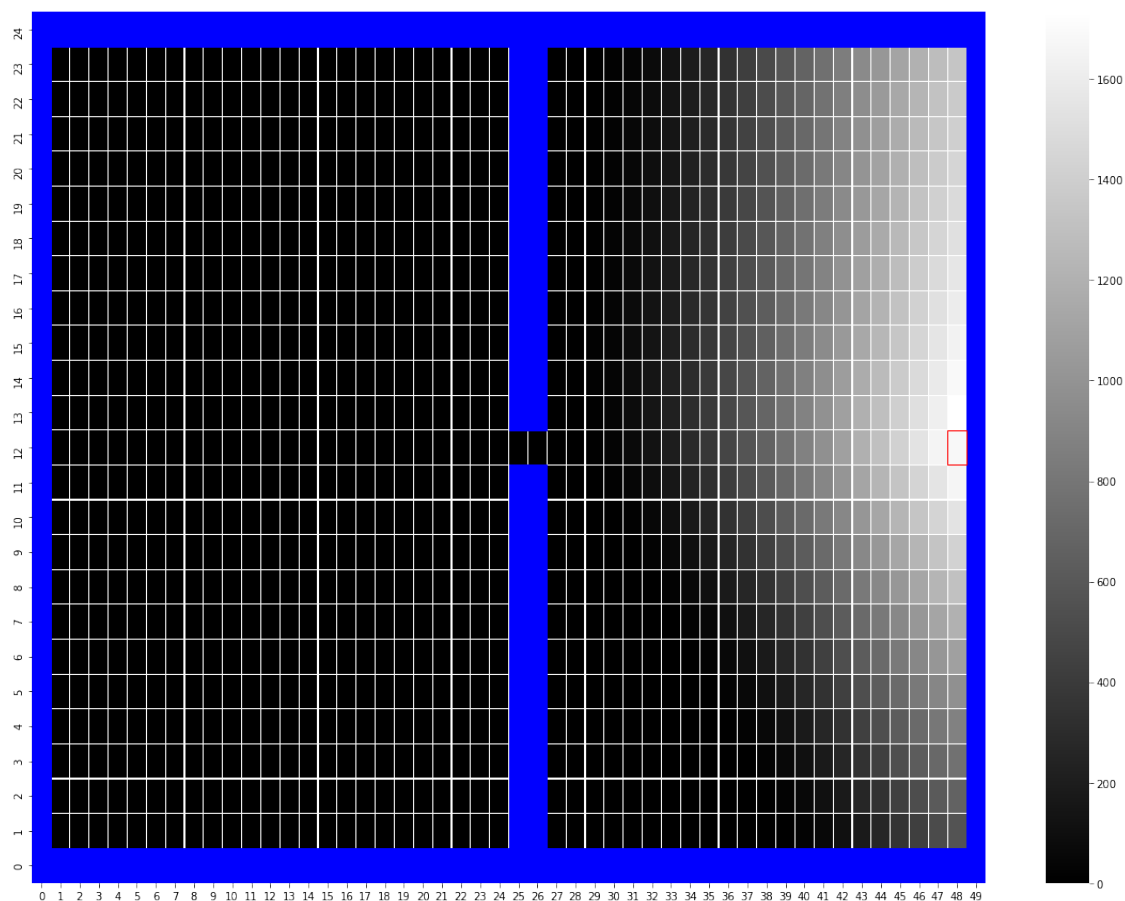


Figure 3: State-Value Function after 20 iterations

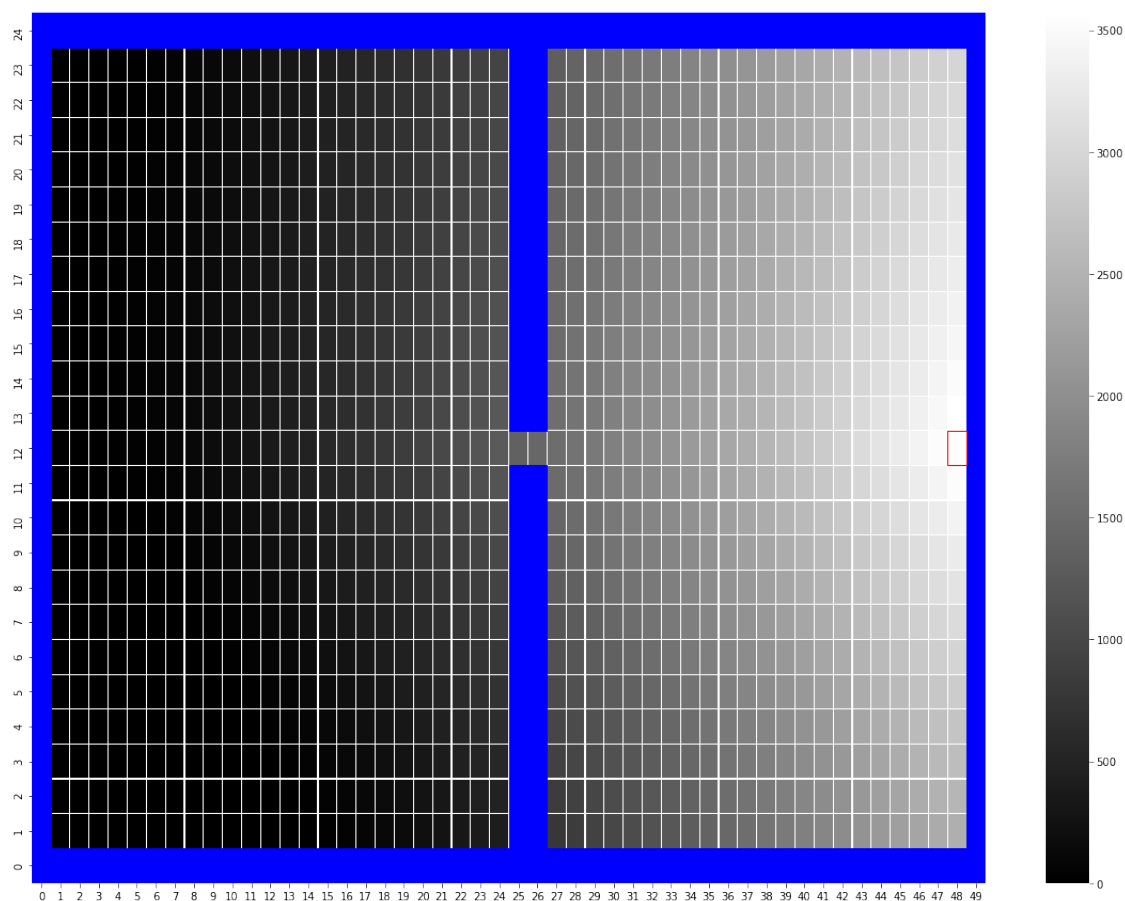


Figure 4: State-Value Function after 50 iterations

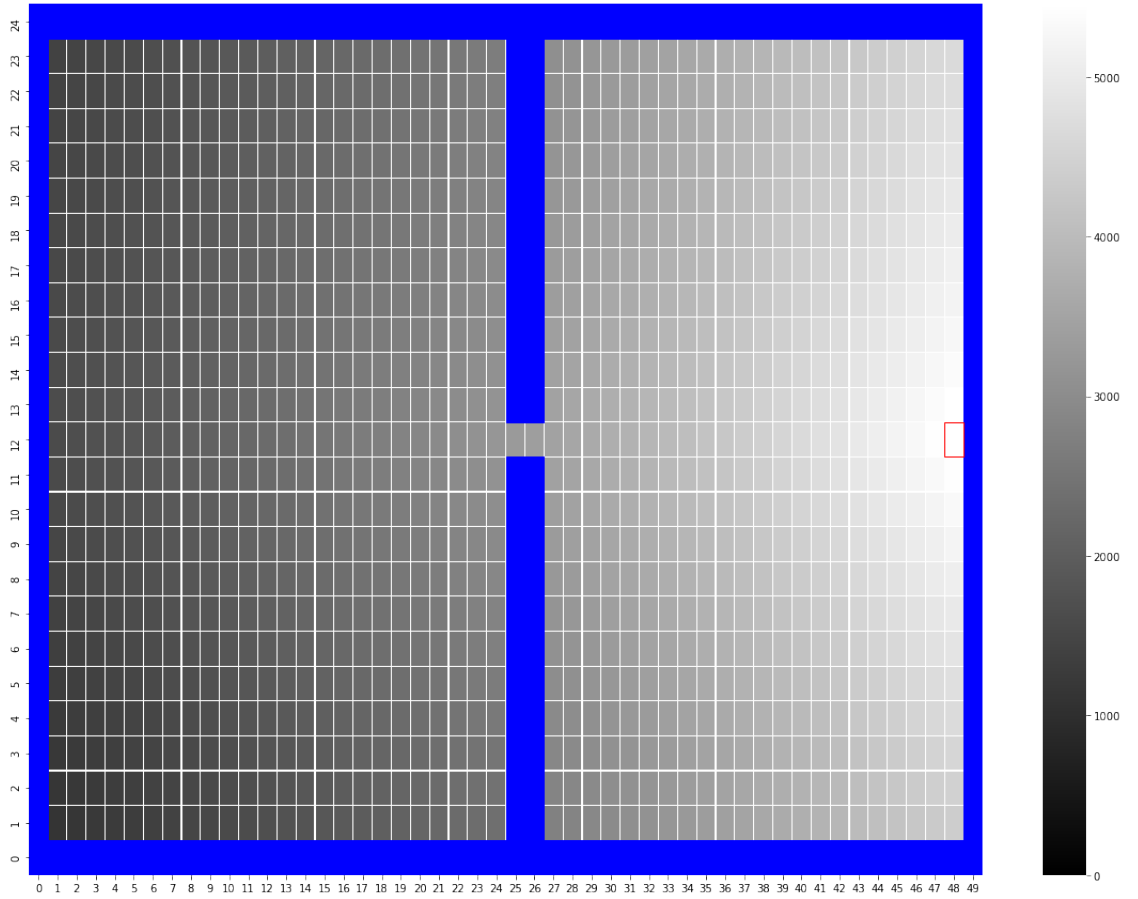


Figure 5: State-Value Function after 100 iterations

1.3 Part c: Policy Simulation

The initial cell in this part is (1, 1).

1.3.1 A Sample Execution of the Policy

A sample execution of the policy obtained after 100 iterations, with $\gamma = 0.99$ is shown in 6

The blue portions indicate the walls and the red cell is the goal.

There are several repeated states and retraced paths in the shown trajectory of the agent, which couldn't be captured.

1.3.2 State-Visitation counts

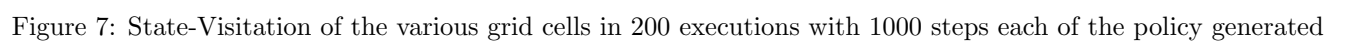
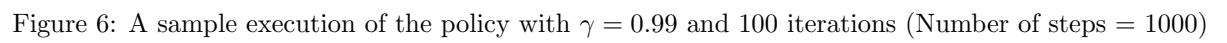
The state-visitation counts for 200 executions (each with 1000 steps) of the policy generated after 100 value-iteration steps with $\gamma = 0.99$ are visualized in the image 7.

To facilitate better visualization of the number of times a state gets visited, the log plot of the state-visitation counts shown in 8, has been generated. It has been generated by adding 1 to each of the state-visitation counts, to make the log function defined.

The blue patches indicate walls and the red cell is the goal cell. The black cells have 0 visitation count.

1.3.3 Inference

1. The Log-plot of the State-visitation count matrix indicates the regions of the grid where the agent has less tendency to traverse.
2. The plot indicates the agent has maximum tendency to stay at the goal state. This is because the goal state has a high positive reward and is not a terminating state. Thus, the visitation of the goal state keeps on rising till the end of 1000 steps.



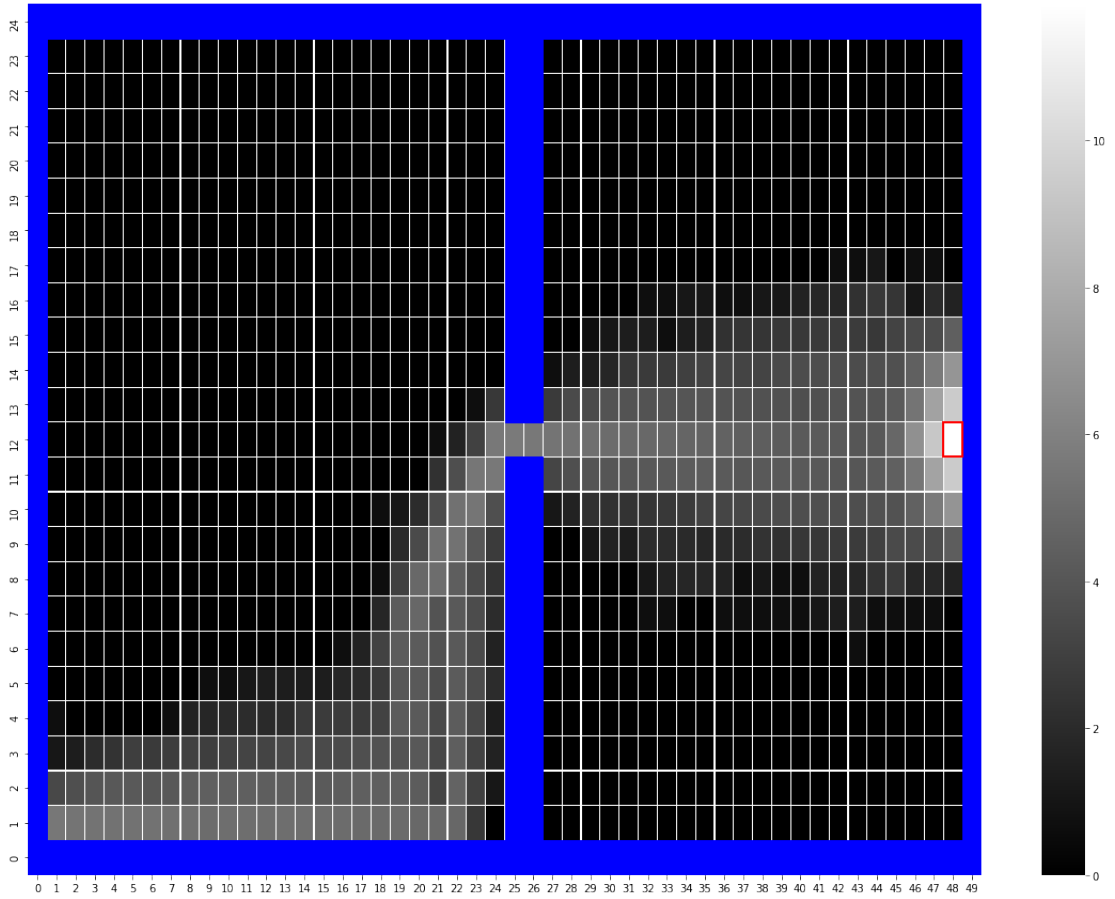


Figure 8: Log plot of state-visitation matrix 7

1.4 Part d: Analysing Value Function and Policy convergence

Randomly sampled value iterations are shown in the table below, for the given parameters.

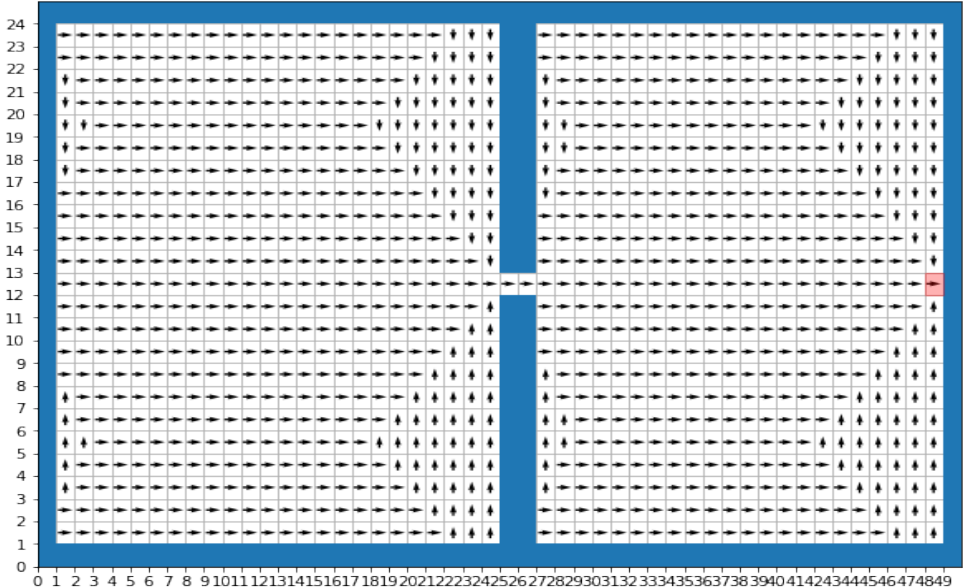
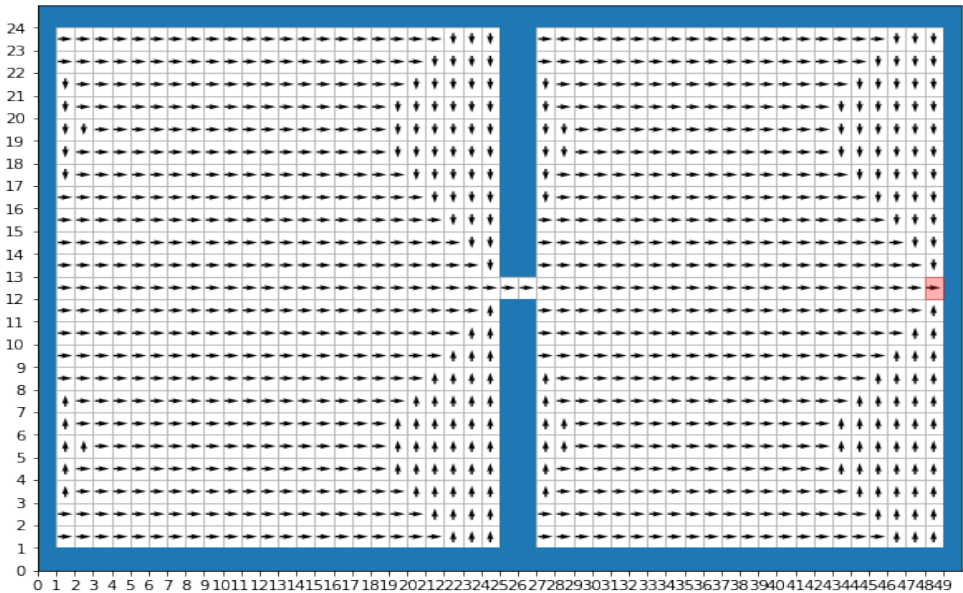
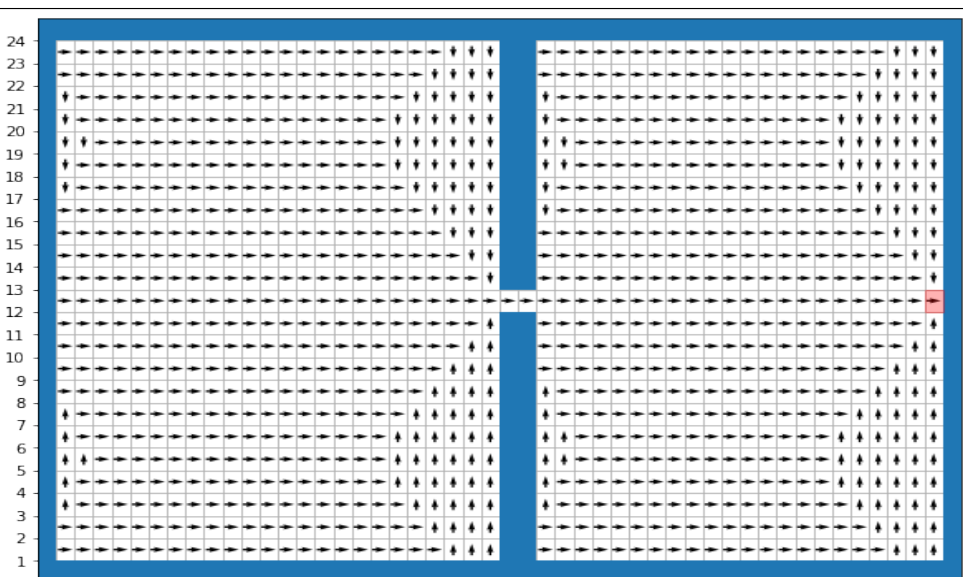
1.4.1 Observations

- $\gamma = 0.99$
- Threshold, $\theta = 0.1$

Iteration	Maxnorm	Policy
1	155.8	

45	56.7	
98	30.2	
101	29.1	

154	15.1	
204	8.1	
271	3.55	

290	2.8	
310	2.2	
370	1.0	

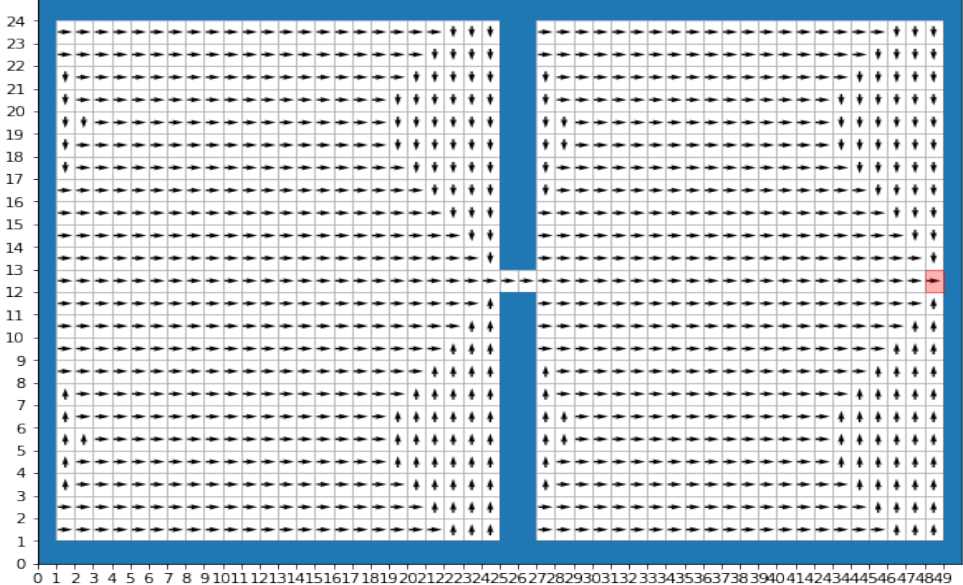
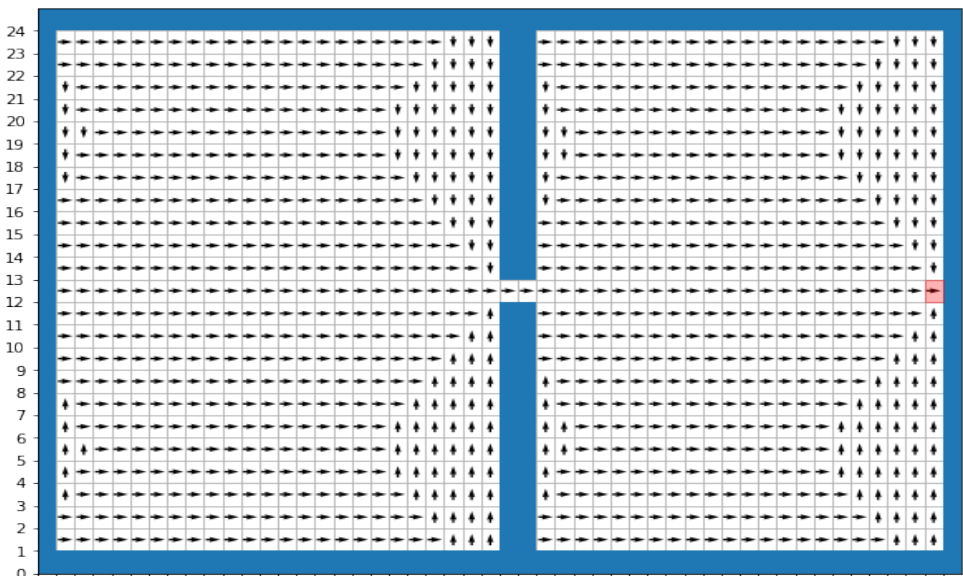
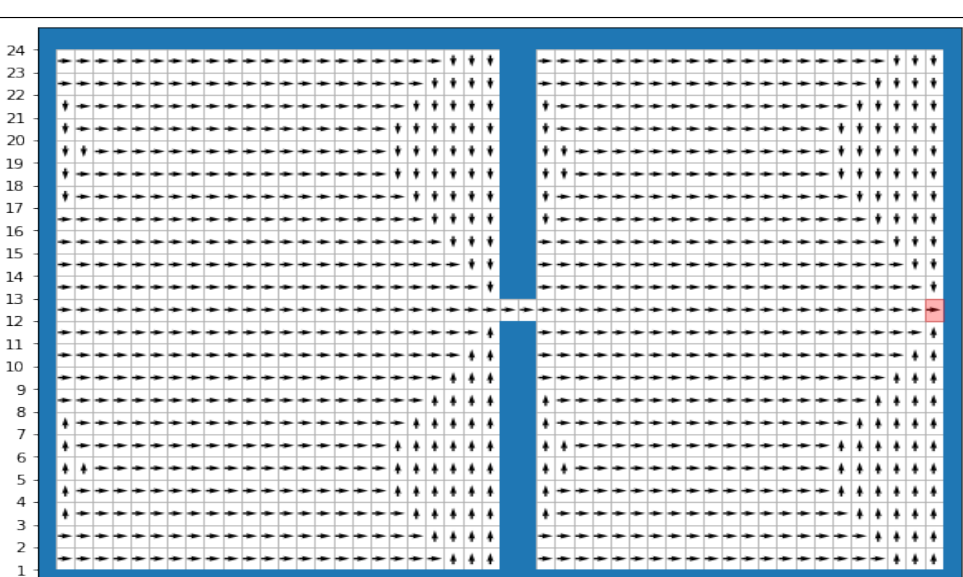
422	0.5	
487	0.24	
560	0.09	

Table 1: Variation of MaxNorm and Policy wrt iteration count for $\gamma = 0.99$

The Policy shown in 1, converges to the final policy after 271 iterations. The maxnorm converges after 560 iterations.

- $\gamma = 0.01$
- Threshold, $\theta = 1e - 50 \approx 0$
- Iterations done = 100

Iteration	Maxnorm	Policy
1	79.57	
5	4.47e-07	
7	3.31e-11	

11	$1.86\text{e-}19$	
14	$1.21\text{e-}25$	
18	$1.54\text{e-}33$	

19	$6.01\text{e-}36$	
20	$7.52\text{e-}37$	
21	$7.34\text{e-}40$	

22	2.86e-42	
23	0	
24	0	

Table 2: Variation of MaxNorm and Policy wrt iteration count for $\gamma = 0.01$

The Policy shown in 2 converges after 19 iterations of Value iteration procedure. The maxnorm becomes a very small value, that it gets approximated to 0.

The variation of the maxnorm with the number of iterations for $\gamma = 0.99$ and $\gamma = 0.01$ is shown in 9

1.4.2 Inference

1. For $\gamma = 0.99$, the maxnorm value falls very slowly as is clear from 9. The gradient in maxnorm vs iterations for $\gamma = 0.01$ is quite high. This is because the state-value function values for $\gamma = 0.01$ is very small for most of the states. This is visible from 1 results. So the maxnorm value will fall quickly for smaller values of γ . As the values are small, maxnorm will also be small. For higher values of γ , the state-values are higher due to higher influence of the goal state reward. So the convergence of the state-value function takes more number of iterations for higher values of γ .
2. As the maxnorm for $\gamma = 0.01$ falls very rapidly, its threshold was kept very low and emphasis was given on the number of iterations.
3. The Policy function is seen to converge after 271 iterations for $\gamma = 0.99$. Policy function converges after 19 iterations for $\gamma = 0.01$. Thus, the policy converges much before the convergence of the value function in both the cases, as expected.
4. As the maxnorm for $\gamma = 0.01$ falls much faster, policy and state-value convergence is much faster as compared to those for $\gamma = 0.99$.

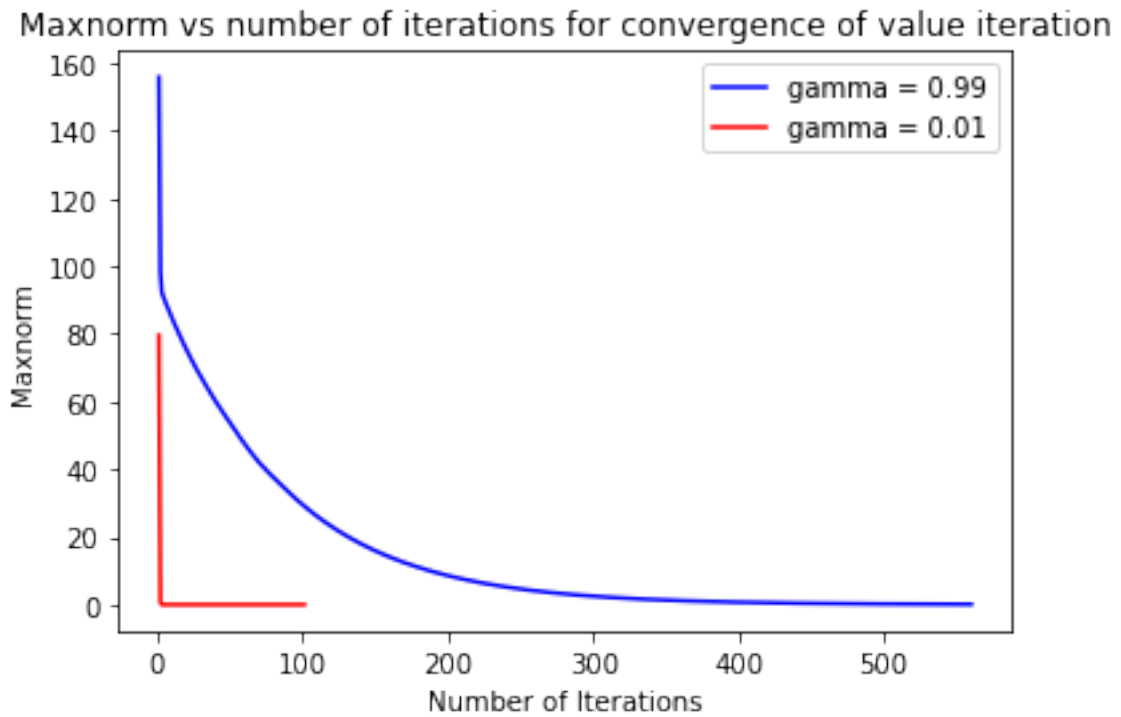


Figure 9: Variation of maxnorm with iterations

2 Q-Learning

2.1 Implementation Details

2.1.1 Background

For updates on each iteration, we have used:

$$\text{qsa_grid}[\mathbf{x}][\mathbf{y}][\text{action}] += \alpha * (\mathbf{r} + \gamma * \text{np.max}(\text{qsa_grid}[\mathbf{x}\mathbf{f}][\mathbf{y}\mathbf{f}]) - \text{qsa_grid}[\mathbf{x}][\mathbf{y}][\text{action}])$$

where:

- α denotes the learning rate
- γ denotes the discount factor
- `action` denotes the action chosen
- (\mathbf{x}, \mathbf{y}) & $(\mathbf{x}\mathbf{f}, \mathbf{y}\mathbf{f})$ denote the states before and after the action
- \mathbf{r} denotes the instantaneous reward obtained

In this case, the learned action-value function, Q , directly approximates q_* , the optimal action-value function, independent of the policy being followed. This dramatically simplifies the analysis of the algorithm and enabled early convergence proofs. The policy still has an effect in that it determines which state-action pairs are visited and updated.

2.1.2 Model

Action Model: Let the intended action be a

$$p(a) = 1 - \epsilon + \epsilon/|A|$$
$$p(A - a) = \epsilon * (|A| - 1)/|A|, \text{ where } A \text{ is the action space.}$$

Reward Model:

$$\begin{aligned} R(\text{Collide with a wall}) &= -1 \\ R(\text{reach goal state}) &= +100 \\ R(\text{otherwise}) &= 0 \end{aligned}$$

Termination: Each episode terminates after 1000 iterations or when the agent reaches the goal state.

2.2 Result

For

- $\alpha = 0.25$
- $\epsilon = 0.05$ (exploration parameter)
- Number of episodes = 4000
- Max length of each episode = 1000

we get the following state-value function (figure 10) and policy (figure 11).

2.2.1 Inference

- The plot for the State-Value function is showing high values not just for the states very close to the goal state, but a lot of other ‘far’ states as well.
- High discount factor value (0.99) as well as a high reward (+100) cause high influence of the goal state on the states far away from the goal.
- Since the episode terminates when the agent reaches the goal state, the max value for the State-Value function is approx 100 (which is the reward for transition into the goal state).
- The policy obtained from MDP was ‘better’ compared to the one obtained by Q-Learning, as it provided a much shorter path to the goal. Even though the policy obtained from MDP wasn’t the best, because it sometimes doesn’t avoid collisions from the walls, it still provides a higher overall reward, which is not true for the Q-Learning case.

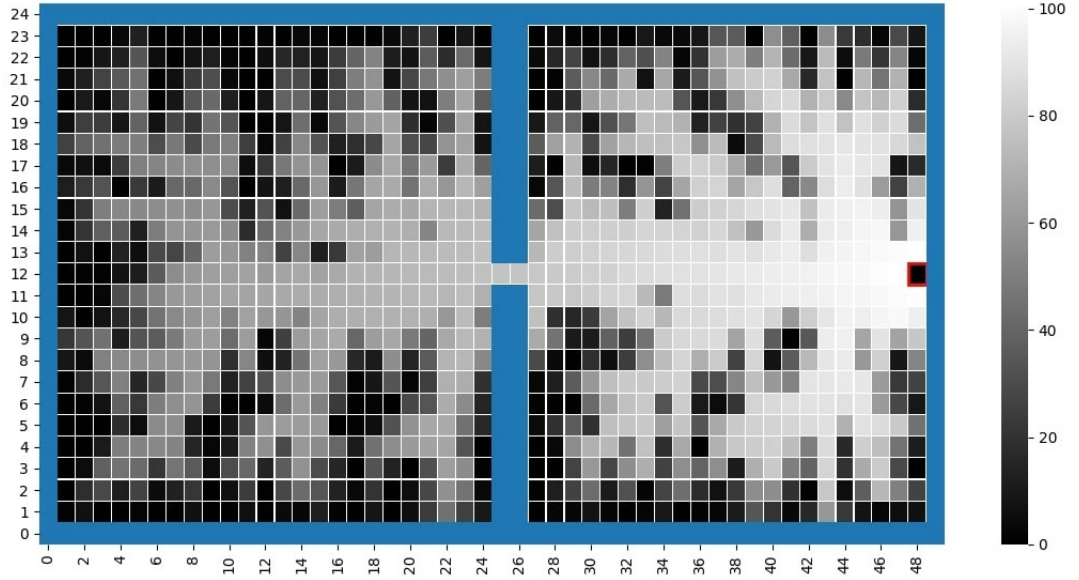


Figure 10: State-Value function for $\epsilon = 0.05$

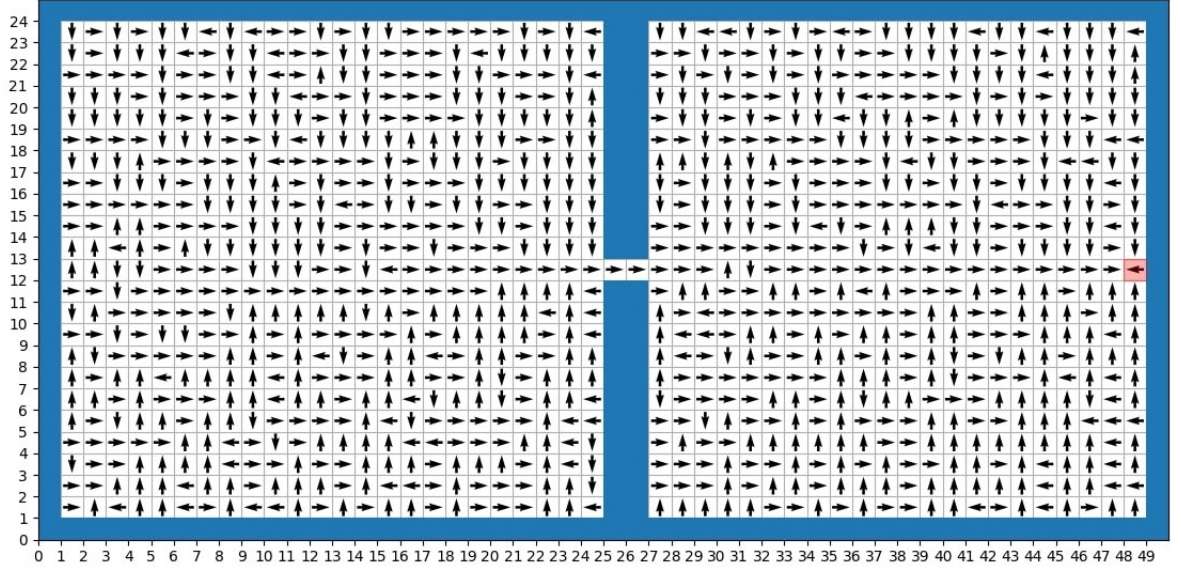


Figure 11: Policy obtained for $\epsilon = 0.05$

2.3 Exploration parameter

Here, we experiment with three different values of the exploration parameter:

- $\epsilon = 0.005$. The state-value function and policy are shown in figure 12 and 13 respectively.
- $\epsilon = 0.05$. The state-value function and policy are shown in figure 14 and 15 respectively.
- $\epsilon = 0.5$. The state-value function and policy are shown in figure 16 and 17 respectively.

2.3.1 Inference

- Exploration parameter is a measure of the probability with which the chosen action will happen. This has been described in section 2.1.2. For the given model, we have 4 actions. Hence, $|A| = 4$.
 - For $\epsilon = 0.005$, the chosen action will happen 99.625% times.
 - For $\epsilon = 0.05$, the chosen action will happen 96.25% times.
 - For $\epsilon = 0.5$, the chosen action will happen 62.5% times.
- For low ϵ , as the agent explores less, it is mostly directed towards the goal, starting from any state. Once it gets the info that goal has good amount of reward, it will direct itself there and get that.

- For very low $\epsilon(0.005)$, the agent doesn't explore enough to get to the goal. It tends to be local. For very high $\epsilon(0.5)$, it explores so much that it loses track of the goal. $\epsilon = 0.05$ serves as a cosy spot for this, providing a good trade-off between exploration and exploitation.
- For $\epsilon = 0.5$, the final policy obtained shows that if the agent starts from the left part of the room ($x < 25$), it can never reach the goal state. We only have two grid points that can be used to cross the centre wall. There would be a lot of negative reward if the agent tries to cross that region, and it instead prefers to stay in the left half itself.

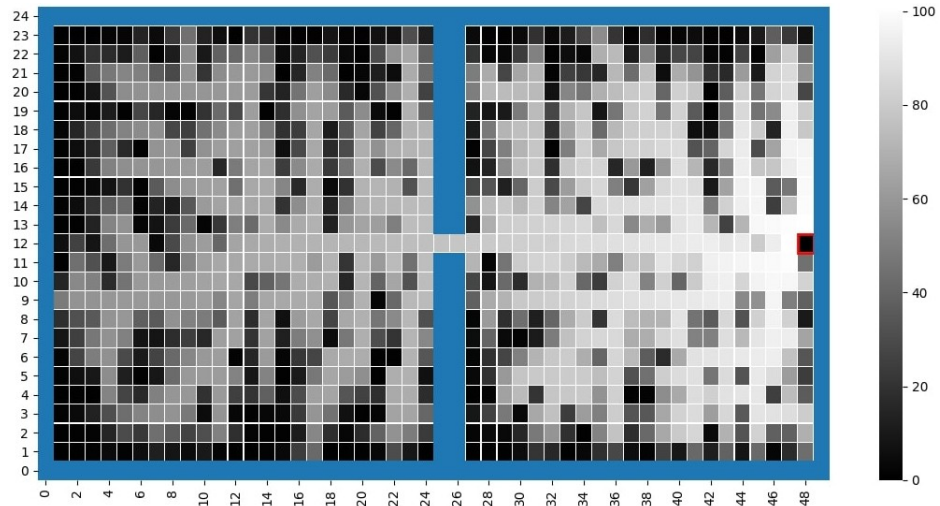


Figure 12: State-Value function for $\epsilon = 0.005$

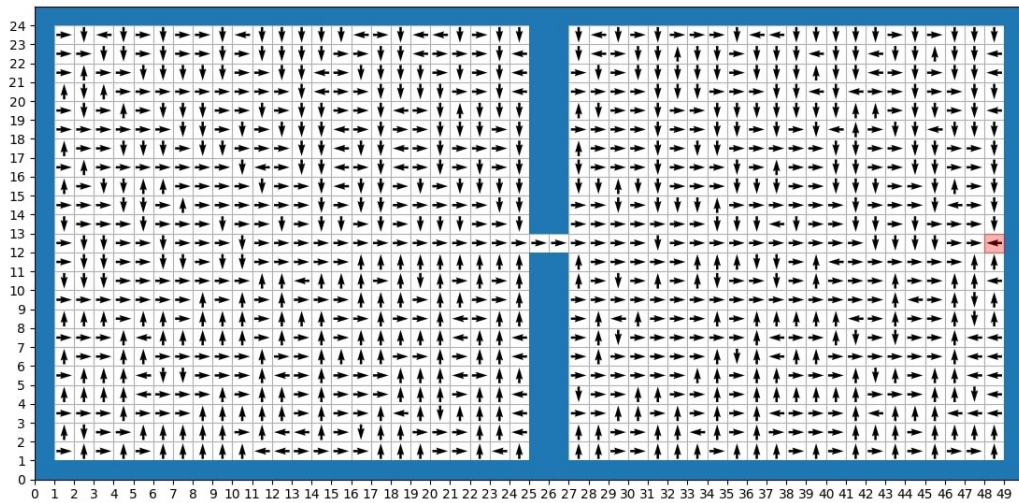


Figure 13: Policy obtained for $\epsilon = 0.005$

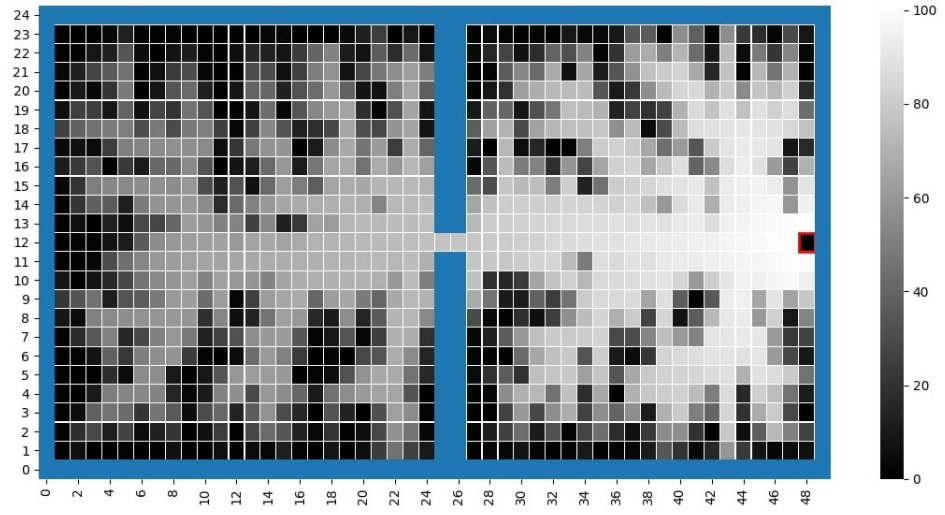


Figure 14: State-Value function for $\epsilon = 0.05$

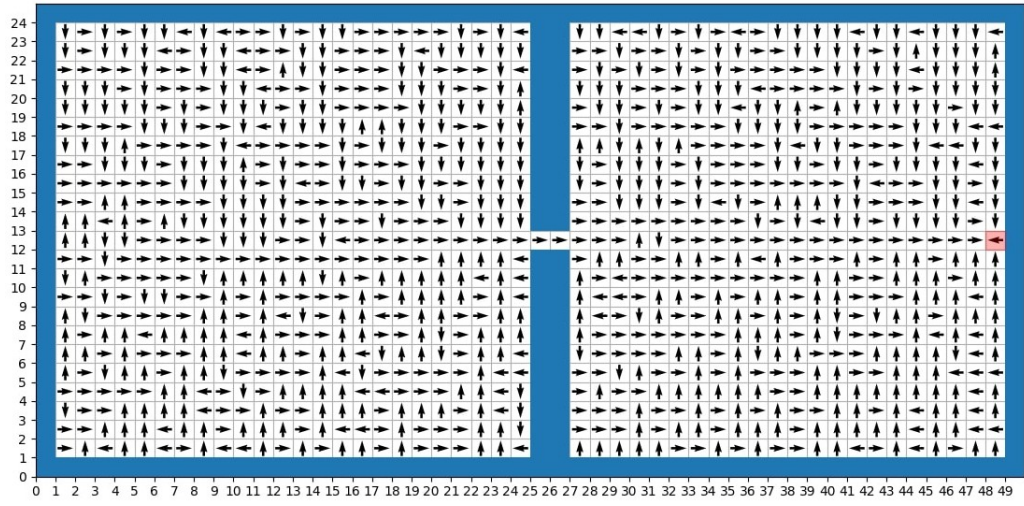


Figure 15: Policy obtained for $\epsilon = 0.05$

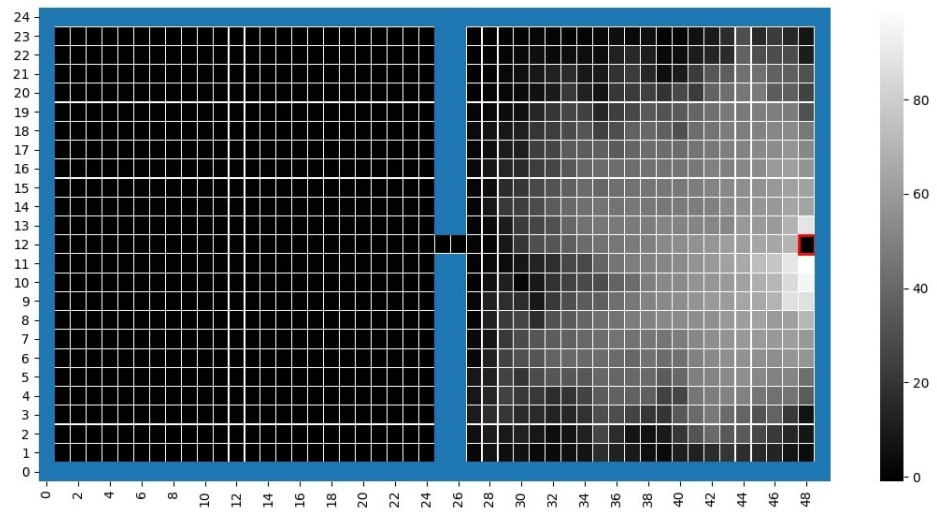


Figure 16: State-Value function for $\epsilon = 0.5$

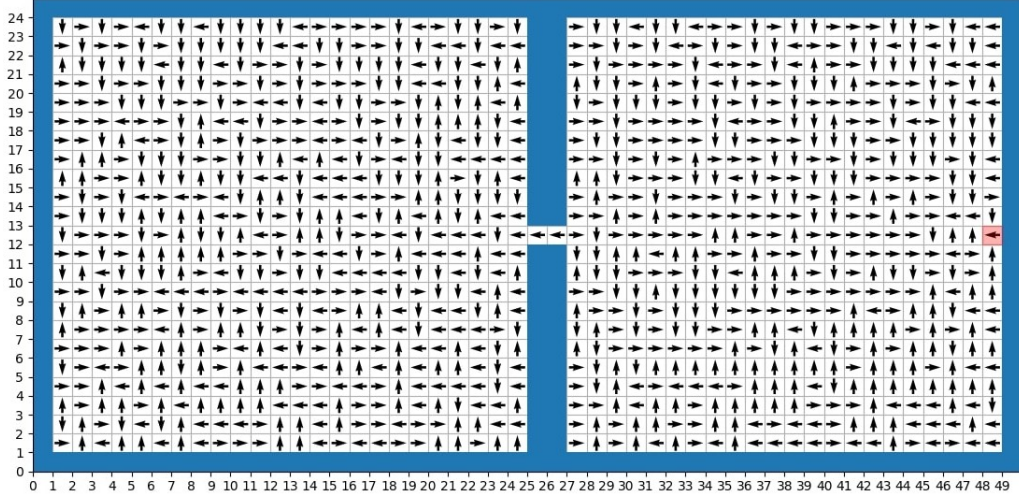


Figure 17: Policy obtained for $\epsilon = 0.5$

2.4 Reward Accumulated

Here, we experiment with two values of the exploration parameter: (note that we have two plots: reward accumulated per episode and average reward accumulated per 10 episodes, and we call them all-reward and avg-reward respectively)

- $\epsilon = 0.05$. The all-reward and avg-reward plots are shown in figure 18 and 19 respectively.
- $\epsilon = 0.5$. The all-reward and avg-reward plots are shown in figure 20 and 21 respectively.

2.4.1 Inference

- As explained in section 2.3.1, for $\epsilon = 0.05$ and 0.5 , the chosen action occurs 96.25% and 62.5% of the times respectively.
- For $\epsilon = 0.05$, it learns the ‘optimal’ path after approximately 800 episodes. Even around 4000th episode, there are instances when the reward accumulated is slightly less than 100 (denoting collisions with the walls), but it is mainly due to the exploratory behaviour.
- For $\epsilon = 0.5$, the arbitrariness in the motion is just too much, and if it starts very close to walls, there can be a lot of collisions with the walls, leading to negative reward. As seen from the policy plot in section 2.3.1, if it starts from the left part of the room ($x < 25$), it can never reach the goal state. Thus, the reward accumulated for those episodes is going to be ≤ 0 . The agent does however obtain +100 reward occasionally when it starts in the right half of the room.

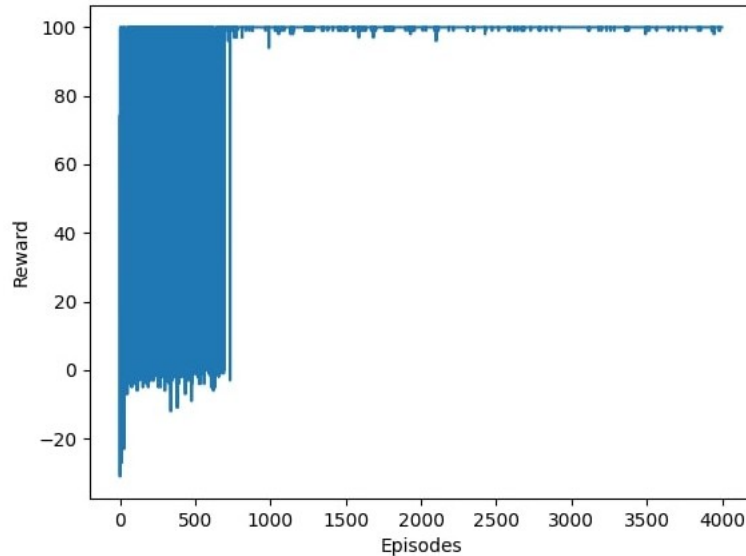


Figure 18: Reward accumulated per episode for $\epsilon = 0.05$

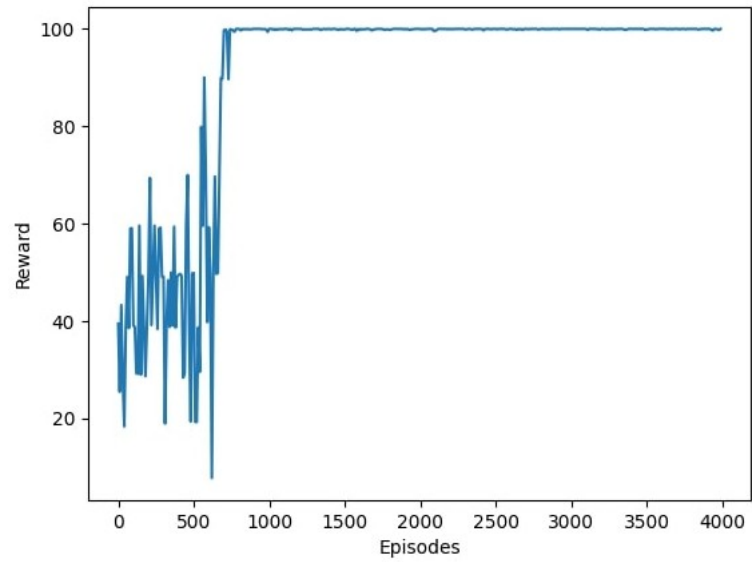


Figure 19: Average reward accumulated per 10 episodes for $\epsilon = 0.05$

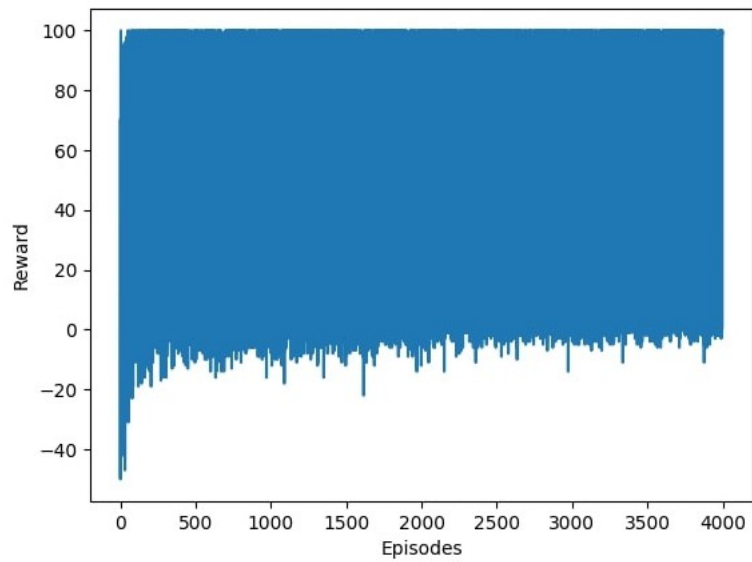


Figure 20: Reward accumulated per episode for $\epsilon = 0.05$

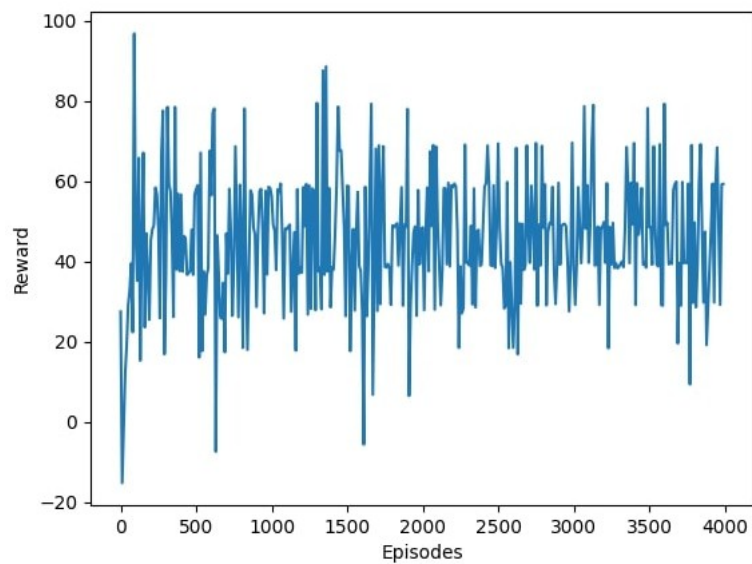


Figure 21: Average reward accumulated per 10 episodes for $\epsilon = 0.05$