

Lab 3

Axel Holmberg (axeho681)

10/4/2020

Environment A

Below are the environments after 10, 100 and 1000 and 10000 runs done with the Greedy Policy. Questions and answers under images.

- What has the agent learned after the first 10 episodes?

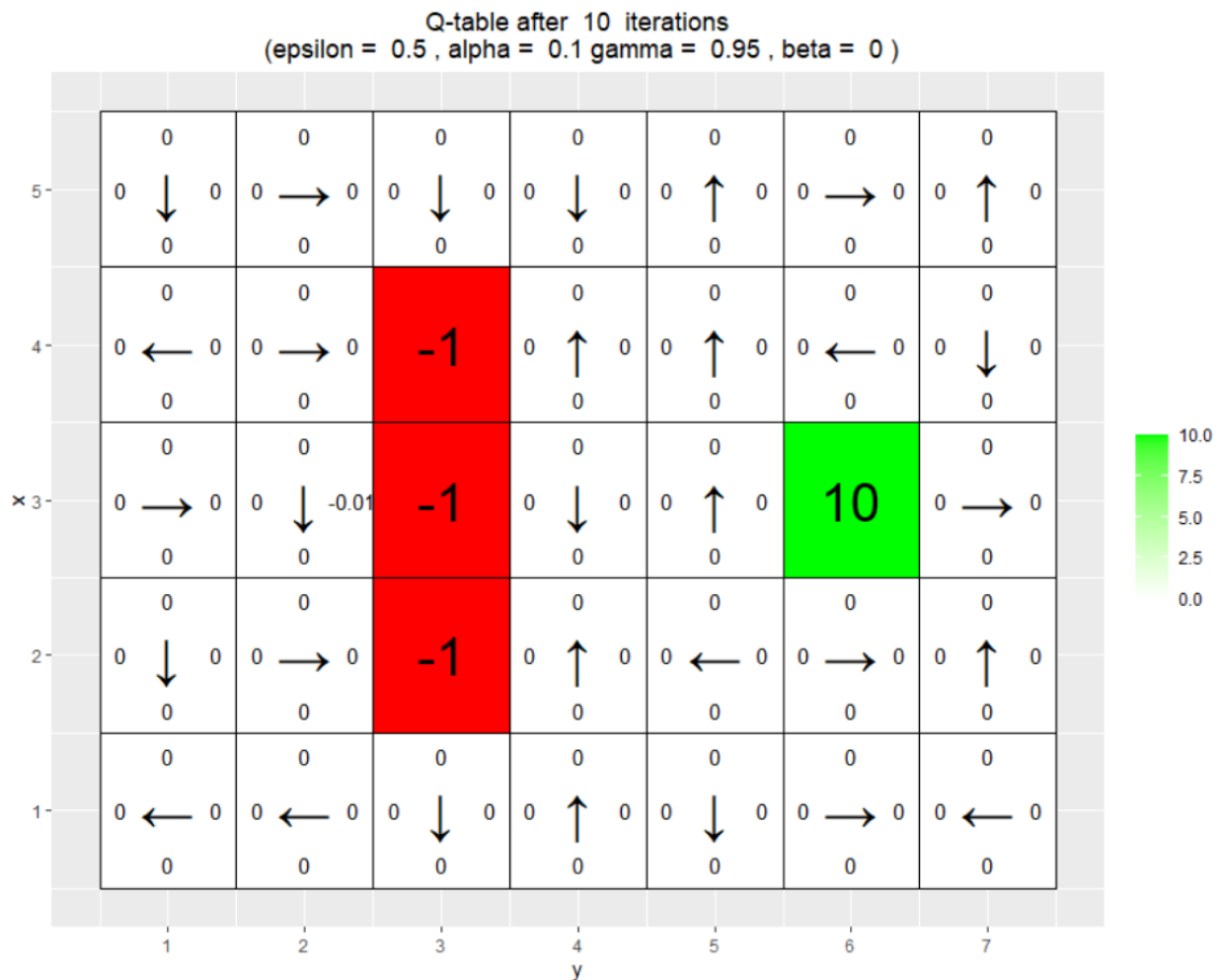


Figure 1: 10 iterations

As one can see in the image above the agent has not learned much in the first 10 iterations. The only q-value that has been updated is after having gotten a negative reward in $x = 3, y = 2$.

- Is the final greedy policy (after 10000 episodes) optimal? Why / Why not ?

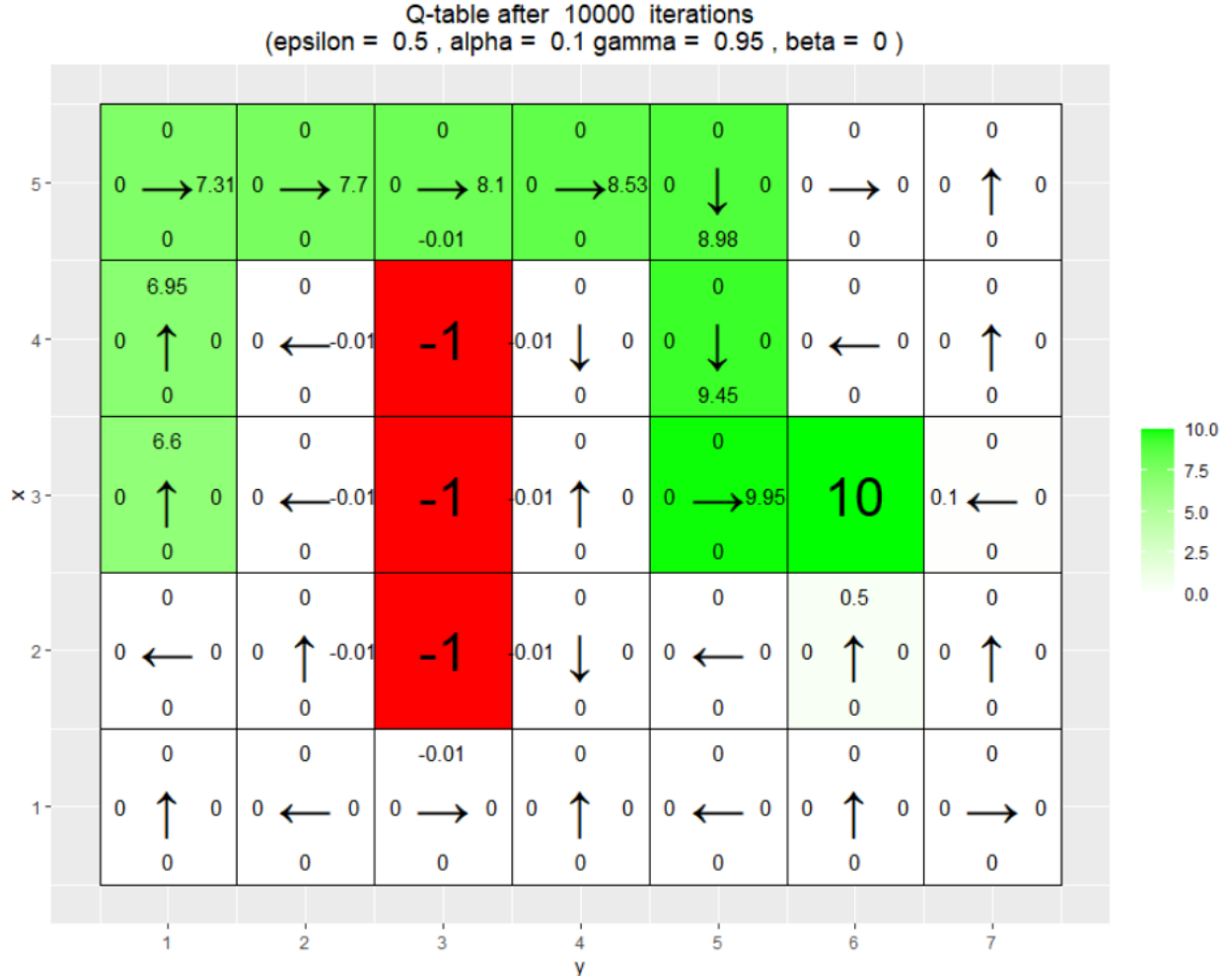


Figure 2: 10 000 iterations

The image above shows the path after 10000 iterations. The agent takes one of the most optimal paths as it takes the least amount of steps possible to get to the reward. Although there are many possible solutions for this where the agent takes exactly 9 steps. Although it would not be optimal if the agent would start in a random tile as it has only optimally done one way to the reward.

- Does the agent learn that there are multiple paths to get to the positive reward ? If not, what could be done to make the agent learn this ?

No, the agent only learns one path. The reason for this is that there is no incentive for it to look for other paths other than the one that it has gotten some reward for as it doesn't have any randomness in the policy for its' movement.

Environment B

As γ denotes the discount factor that means that it adjusts how much of the values of the Q-table should affect the new updated Q-value, called correction. One can see this quite clearly if one compares $\gamma = 0.75$ with $\gamma = 0.95$.

The big difference one can see in the images above is how all the q-values are a lot higher when $\gamma = 0.95$.

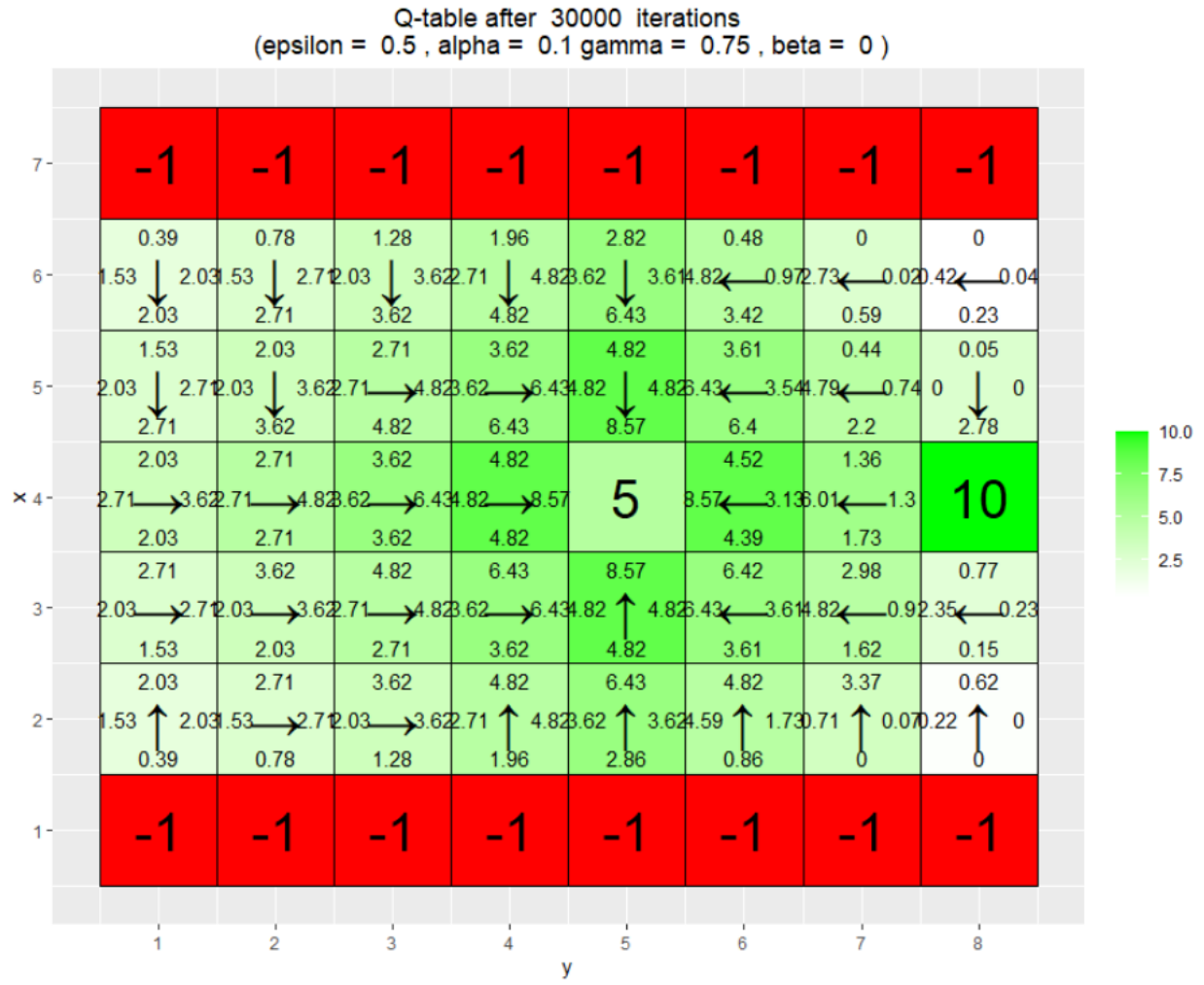


Figure 3: $\gamma = 0.75, \epsilon = 0.5$

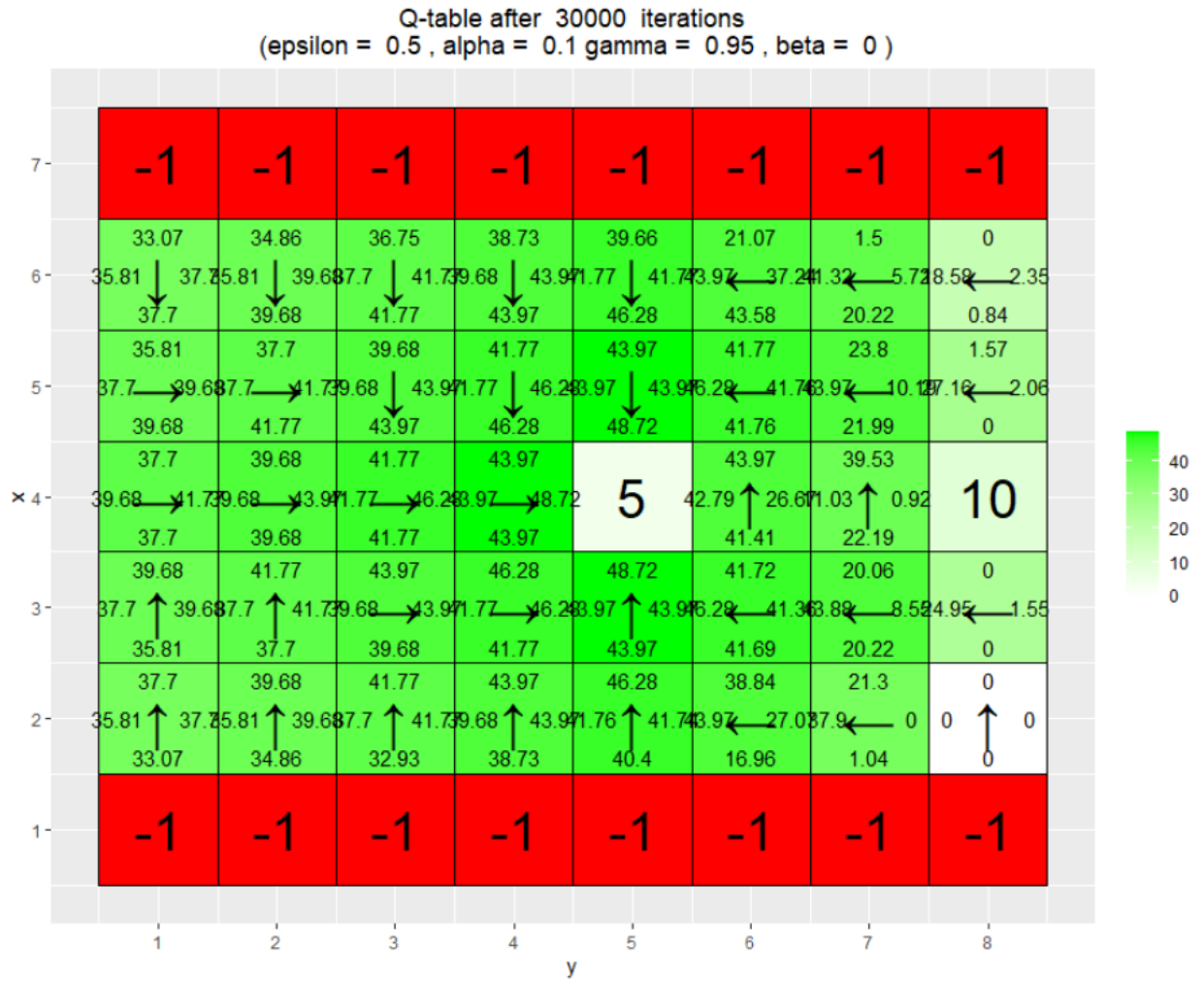


Figure 4: $\gamma = 0.95, \epsilon = 0.5$

One more thing that is clearly visible in these images is the effect of ϵ . ϵ is the threshold for exploration. As $\epsilon = 0.5$ in both of the above that means that in each evaluation of the ϵ -greedy policy there is a 50% chance that it instead takes a random step. That means that it has a high chance of exploring more of the tiles. If one compares this with where $\epsilon = 0.1$ like the images below the results is quite clear.

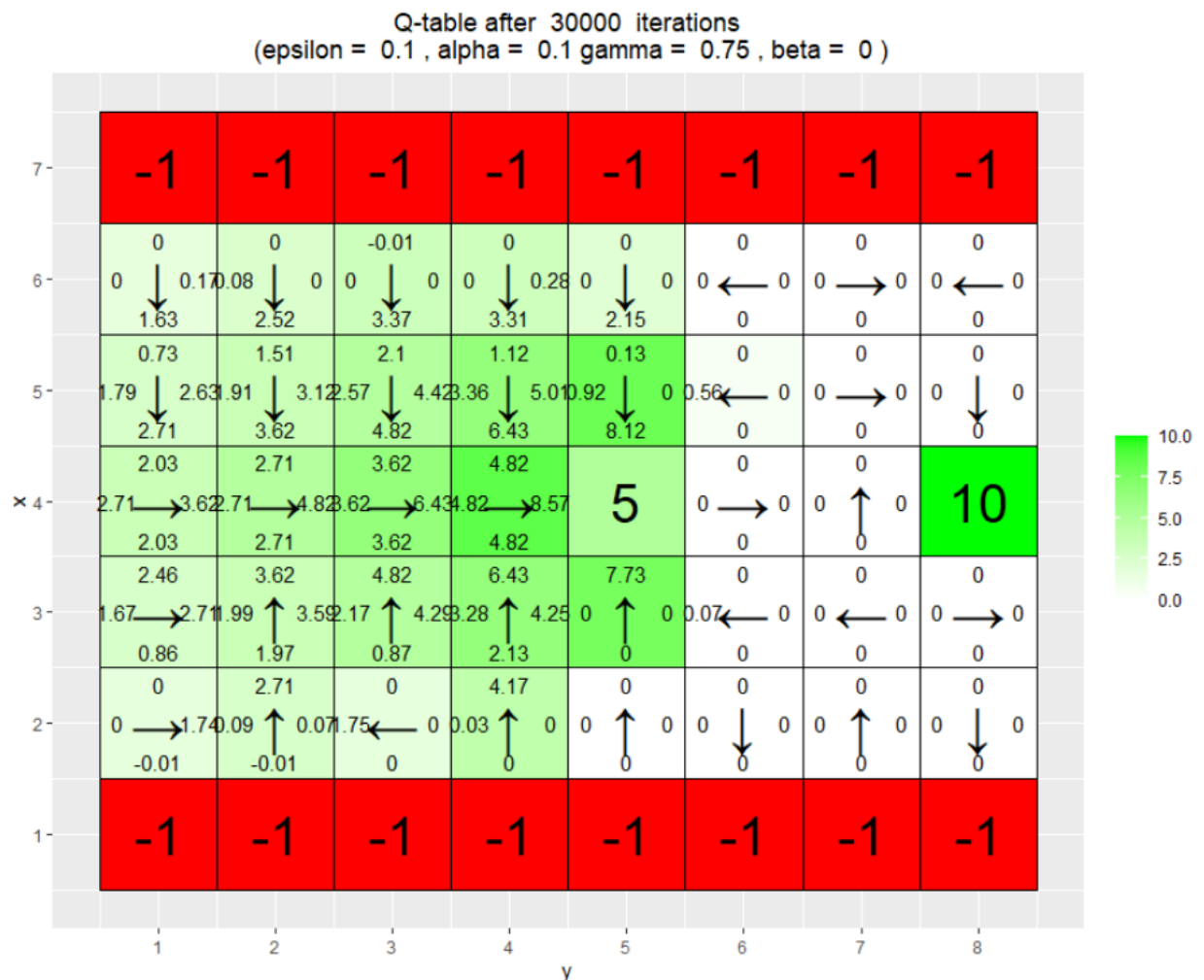


Figure 5: $\gamma = 0.75, \epsilon = 0.1$

As one can see it almost never takes a step beyond the tile with a reward of 5 as the low exploration rate makes it go straight to the tile with a reward of 5. This can also be seen in the graph below showing the rolling mean of the rewards where it almost always 5.

One can also look at the correction graph. The correction graph (see below) shows the correction of each step. This shows how much of a correction that occurs with each step.

The difference between these two are quite significant. The higher value of ϵ leads to more variation and mainly a higher correction overall. What this means is that with each episode the correction of the q values are higher. The reason for is that the higher the probability of acting greedily is the more different paths and more will be discovered and the more correction is done each step and episode.

There are more to be analysed from the data and the graphs, but above is the key takeaways that I did of the graphs.

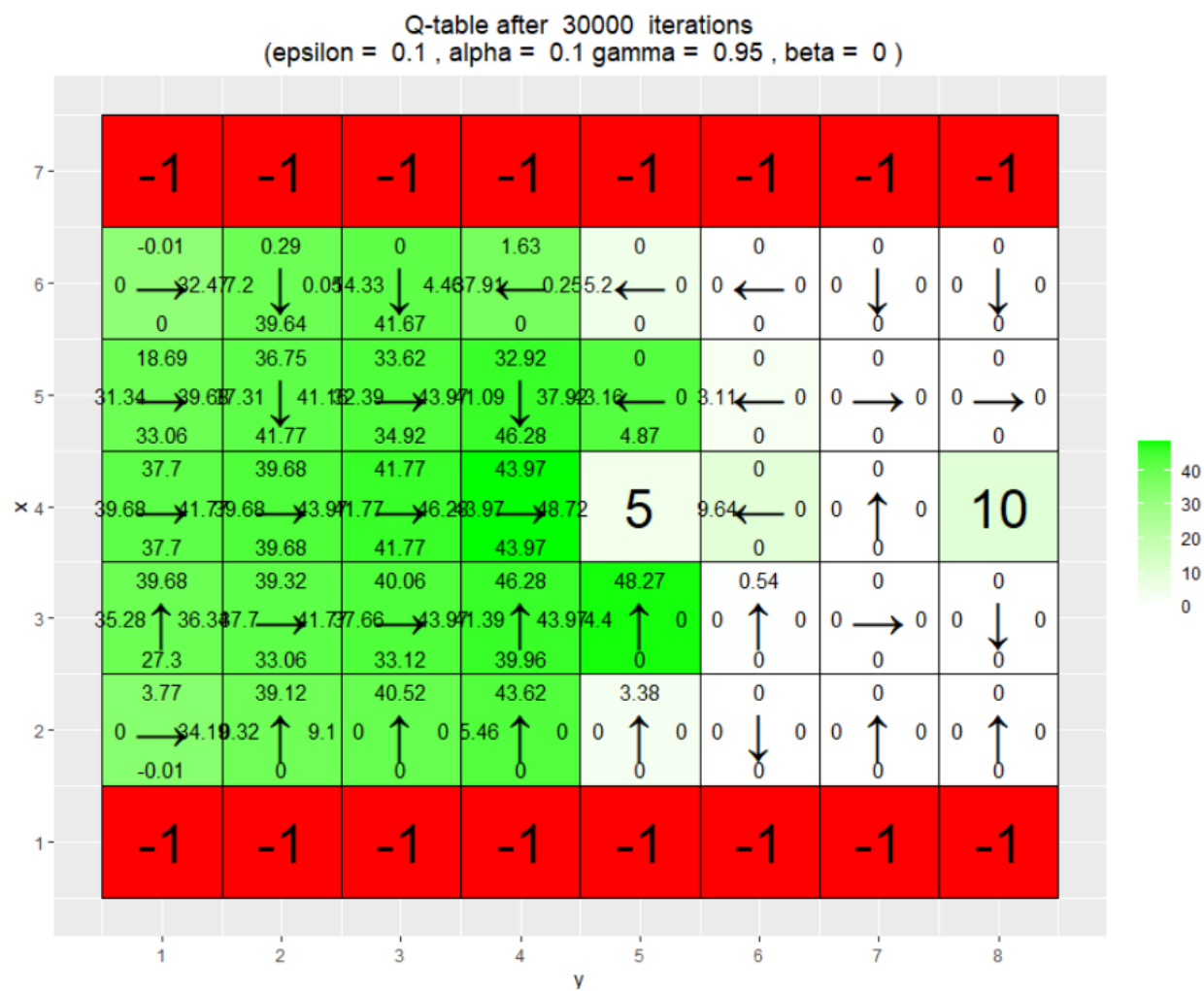


Figure 6: $\gamma = 0.95, \epsilon = 0.1$

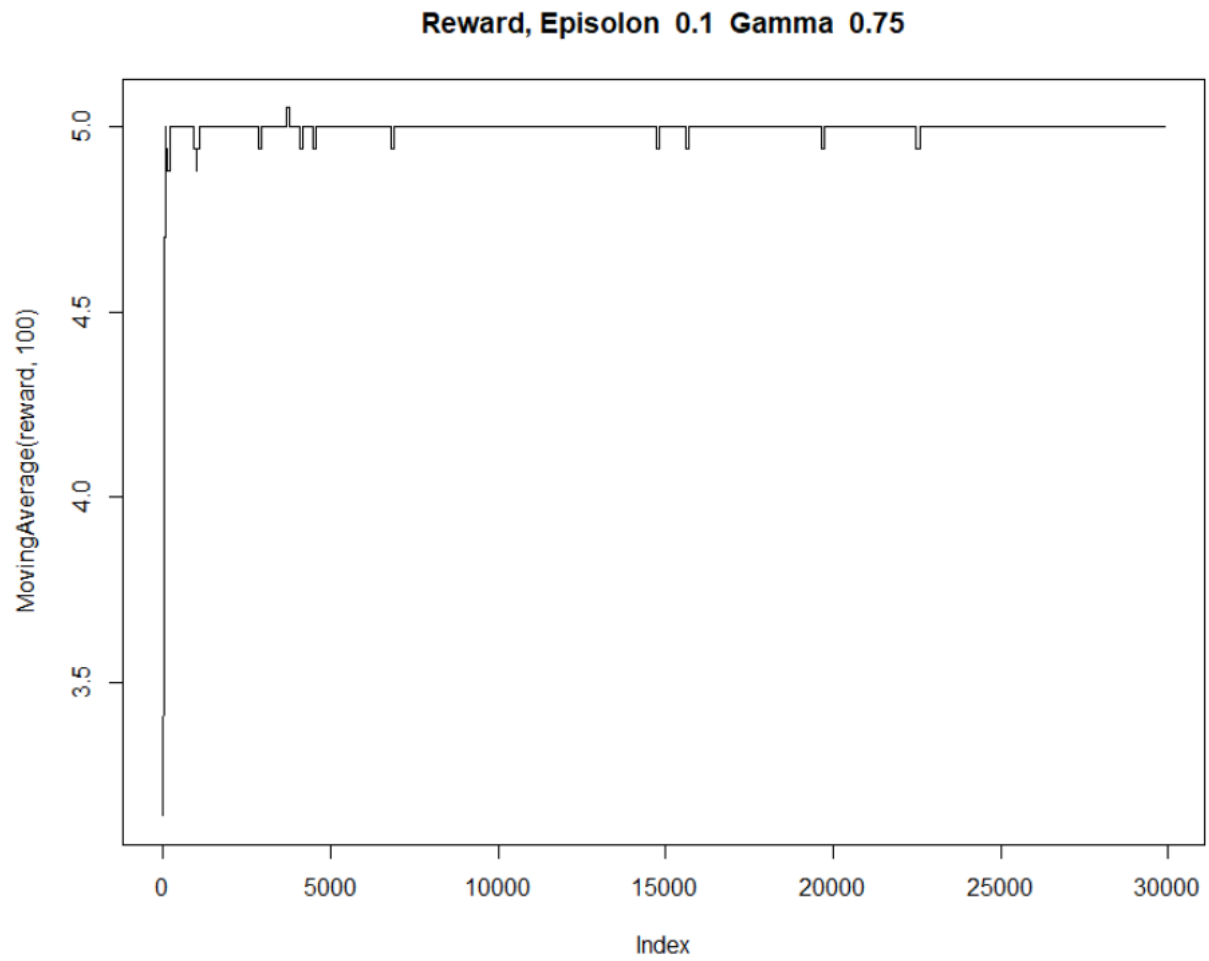


Figure 7: $\gamma = 0.75, \epsilon = 0.1$

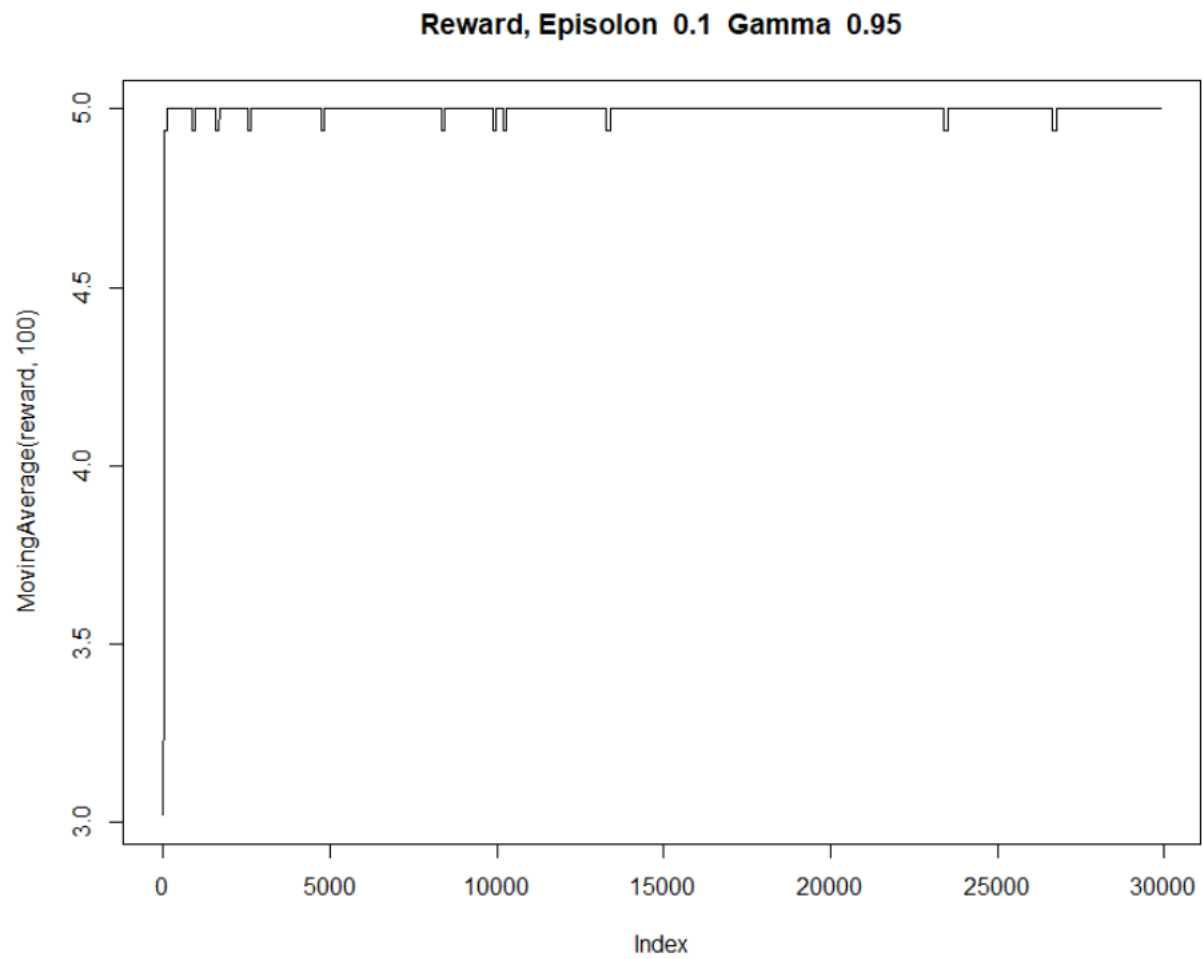


Figure 8: $\gamma = 0.95, \epsilon = 0.1$

Correction, Episolon 0.1 Gamma 0.95

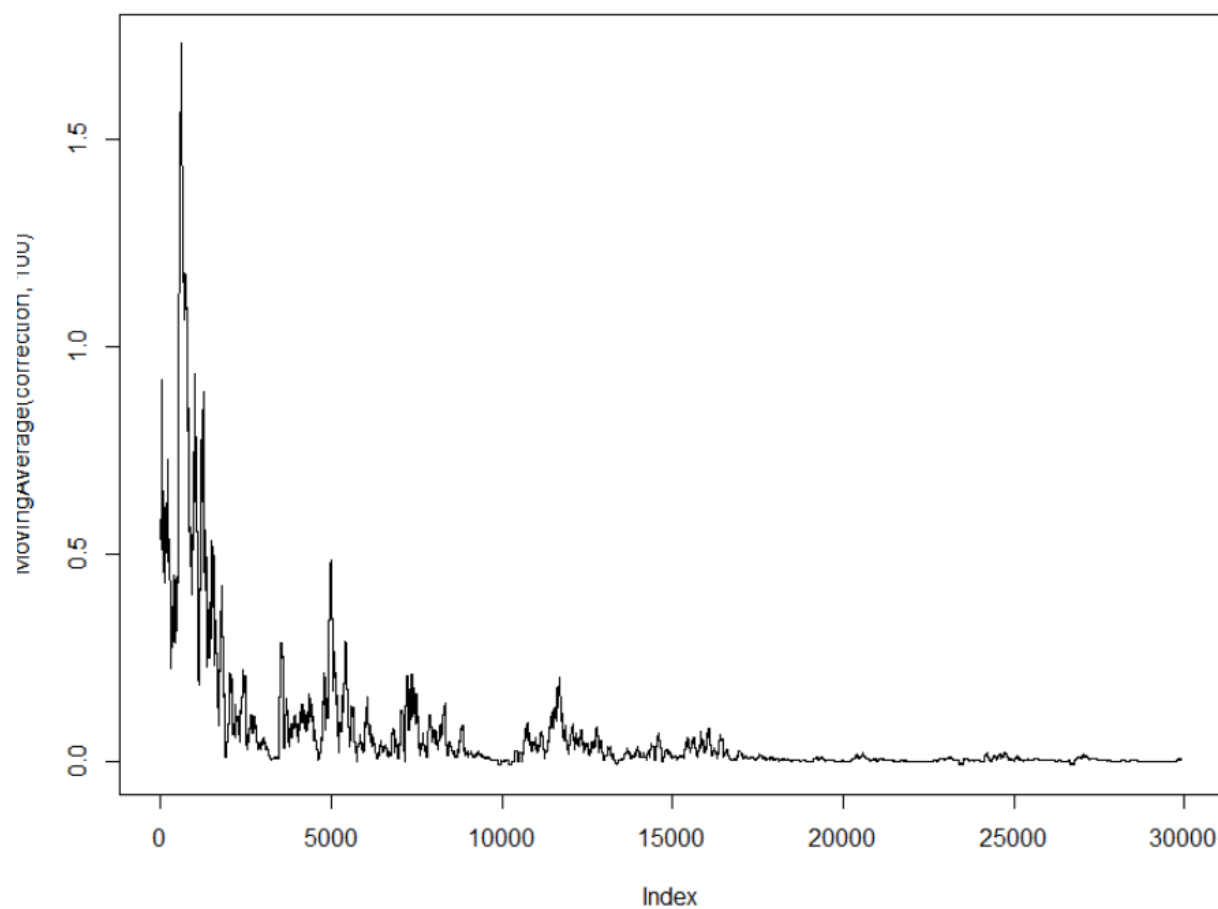


Figure 9: $\gamma = 0.95, \epsilon = 0.1$

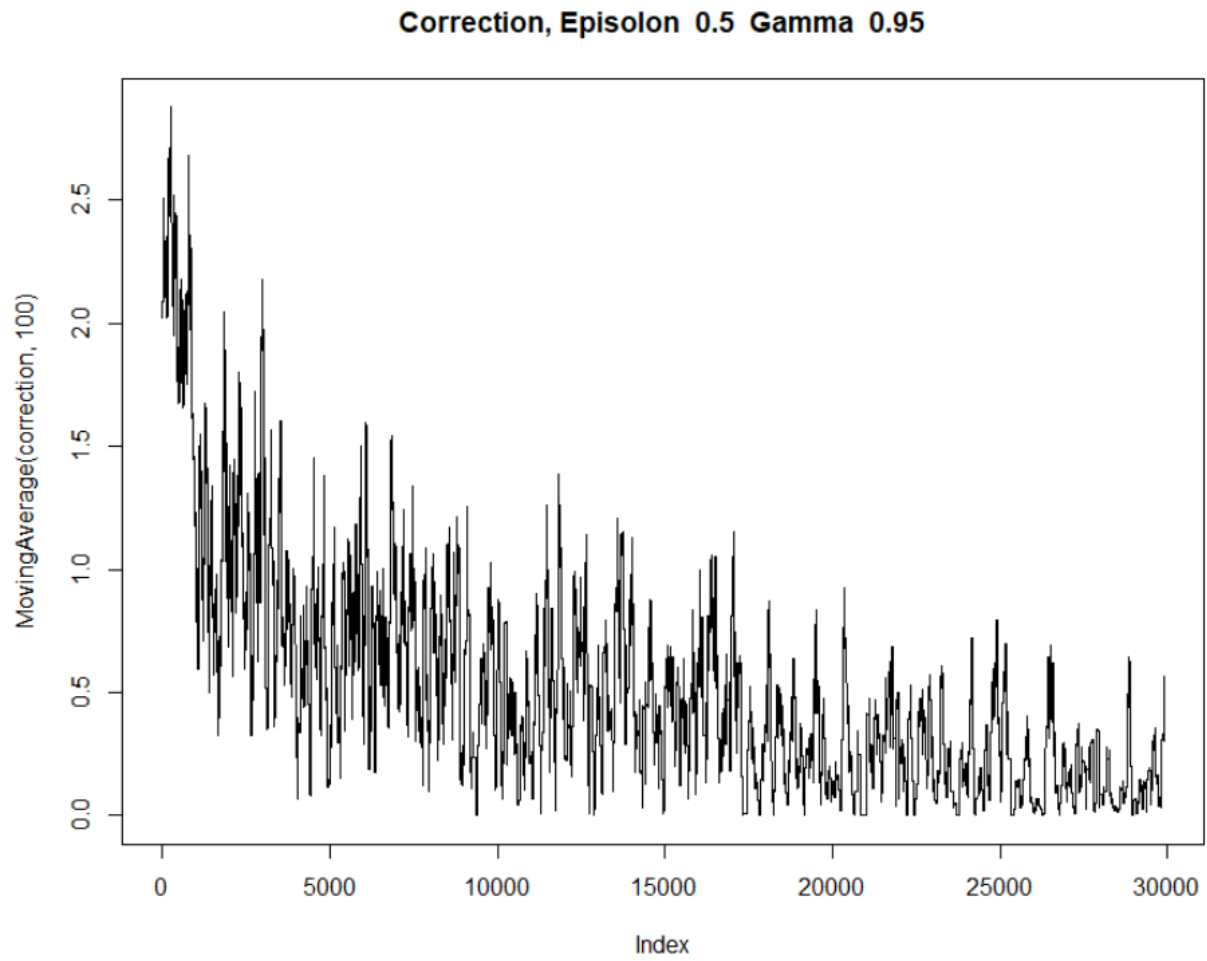
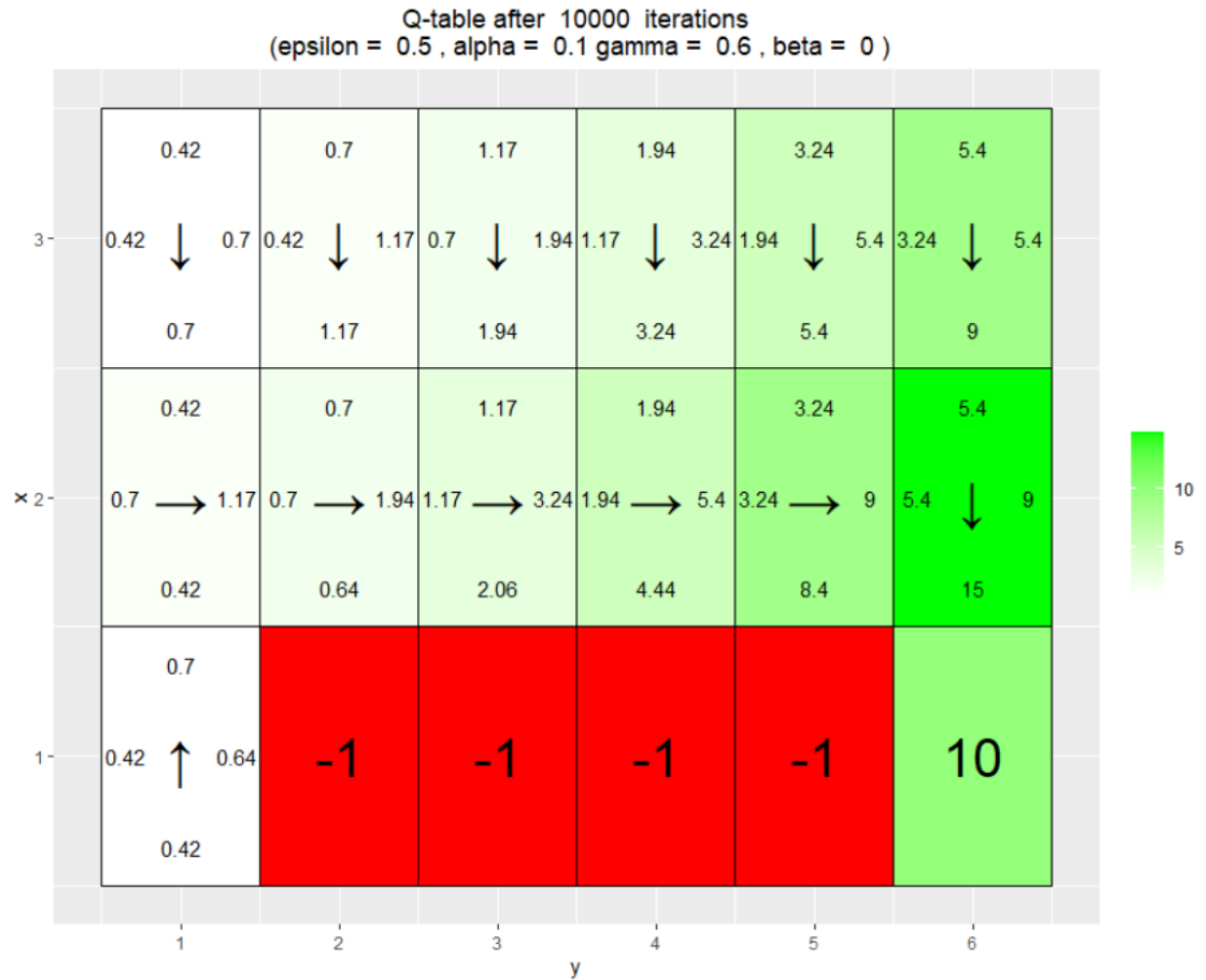


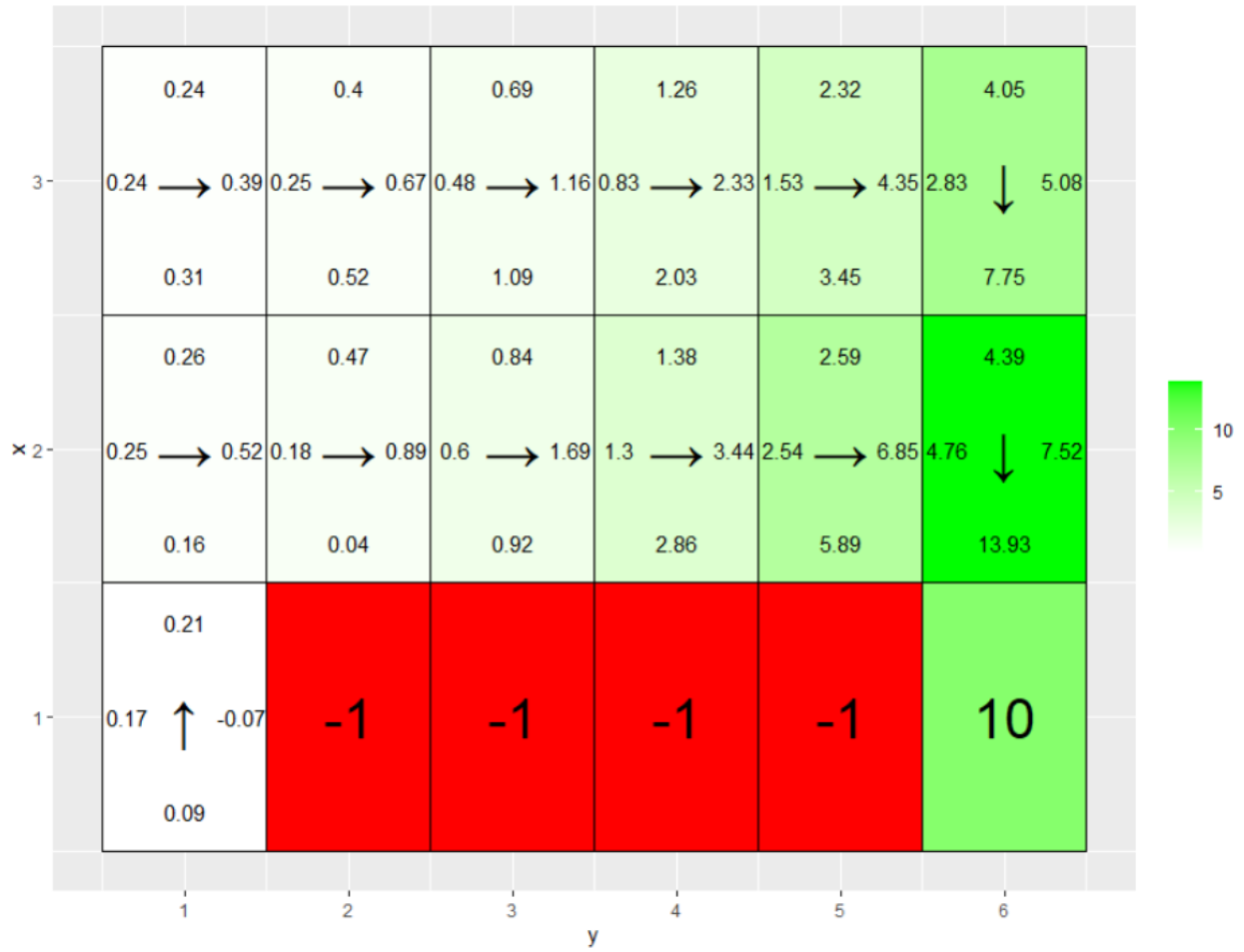
Figure 10: $\gamma = 0.95, \epsilon = 0.5$

Environment C

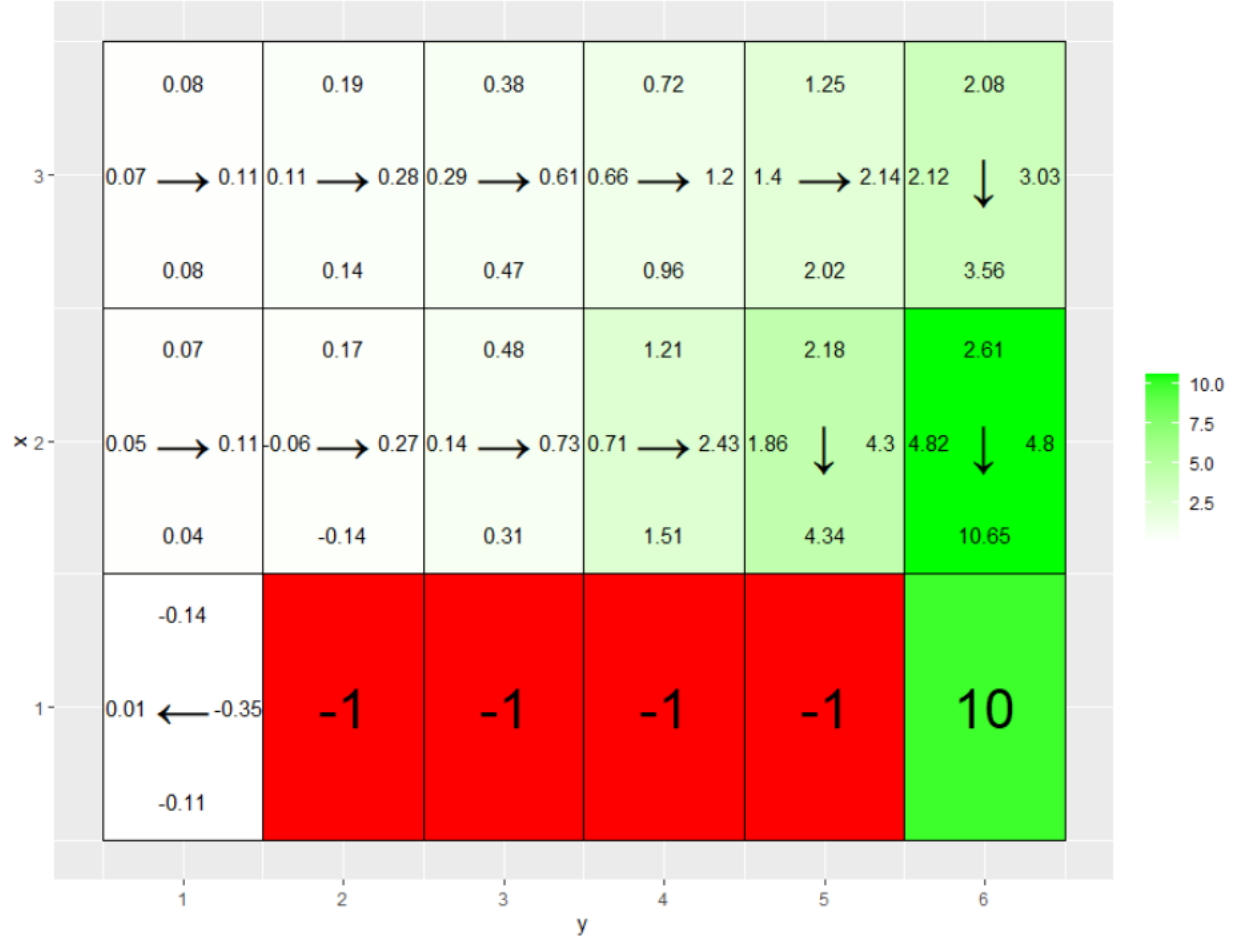
β is the slipping factor. The slipping factor can make the action “slip” and make the action be something different. The β is the probability of the agent slipping to the side when trying to move with each step.

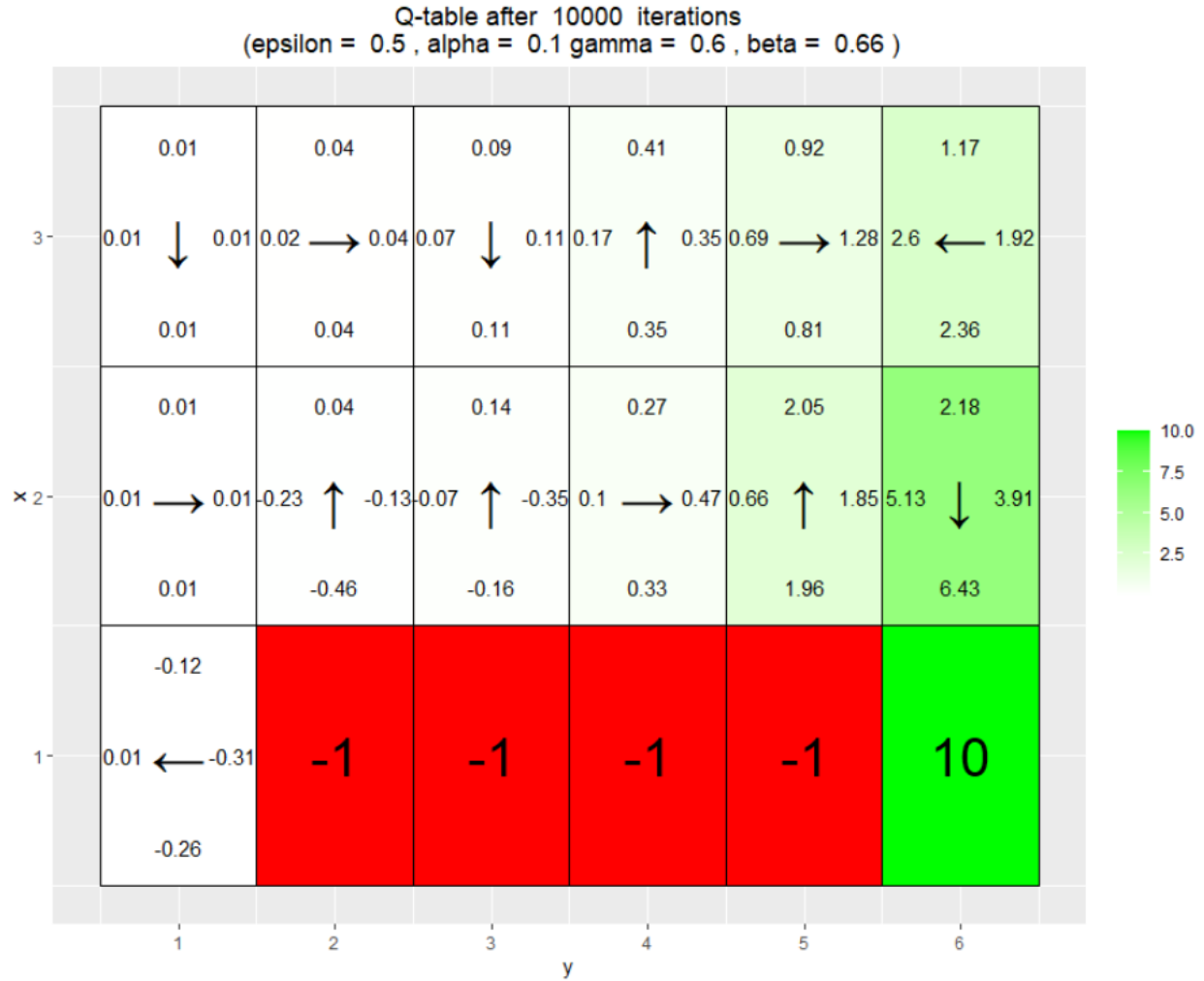


Q-table after 10000 iterations
 (epsilon = 0.5 , alpha = 0.1 gamma = 0.6 , beta = 0.2)



Q-table after 10000 iterations
(epsilon = 0.5 , alpha = 0.1 gamma = 0.6 , beta = 0.4)





Environment D

- *Has the agent learned a good policy? Why / Why not?*

Yes, the agent has learned a good policy as in almost all of the cases the agent would end up at the goal no matter what the starting position would be.

Example:

Example with one bad probability in (1,4):

- *Could you have used the Q-learning algorithm to solve this task?*

TODO

Environment E

- *Has the agent learned a good policy? Why / Why not?*

The policy is worse than in Environment D as there are a lot more cases where if the agent were placed in a random tile in the grid they would not end up at the goal. The reason for this is that .. TODO

Example of a bad policy where the agent would end up almost stuck at the right edge most of the time:

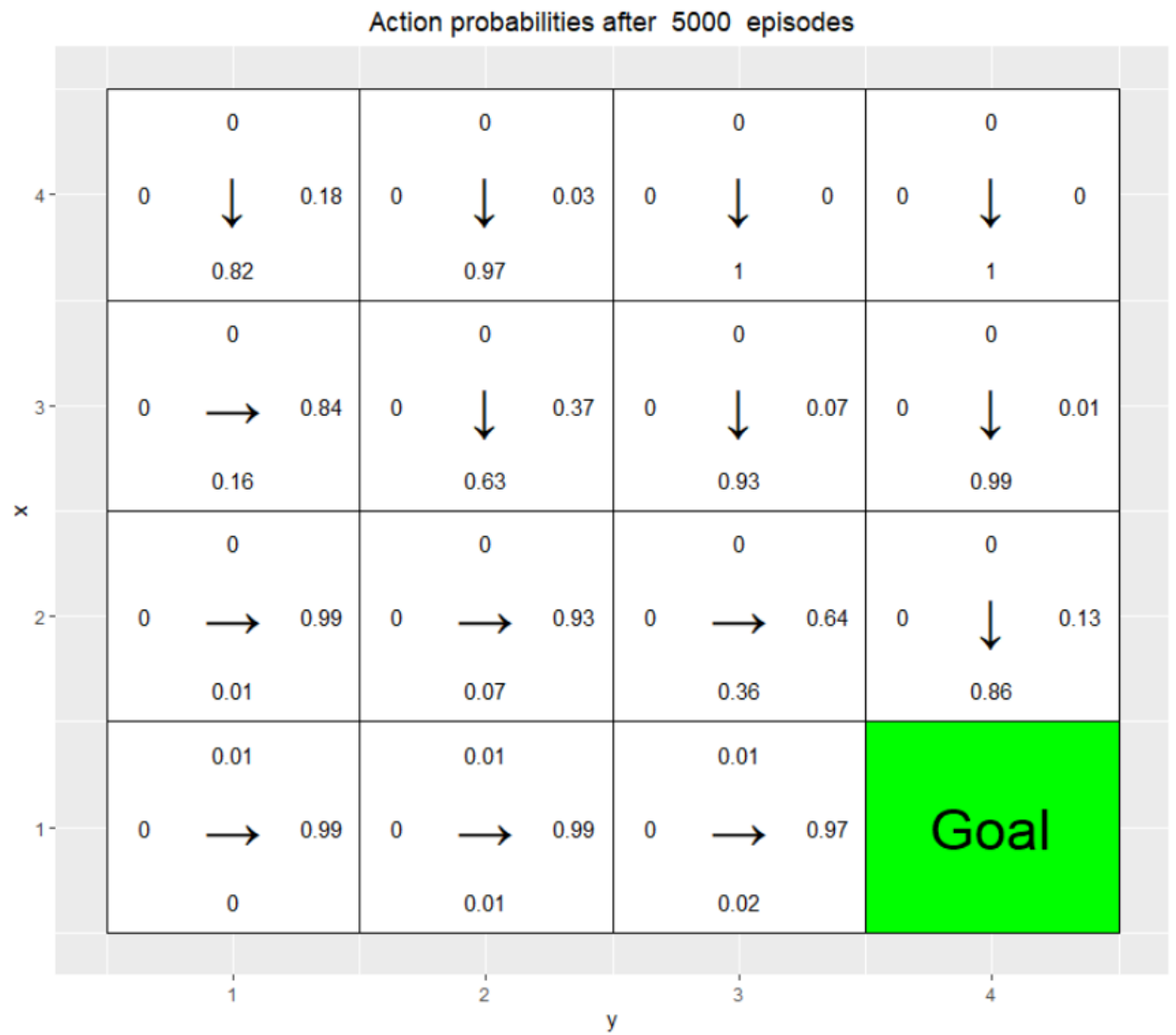


Figure 11: Reinforce with goal in (1,4)

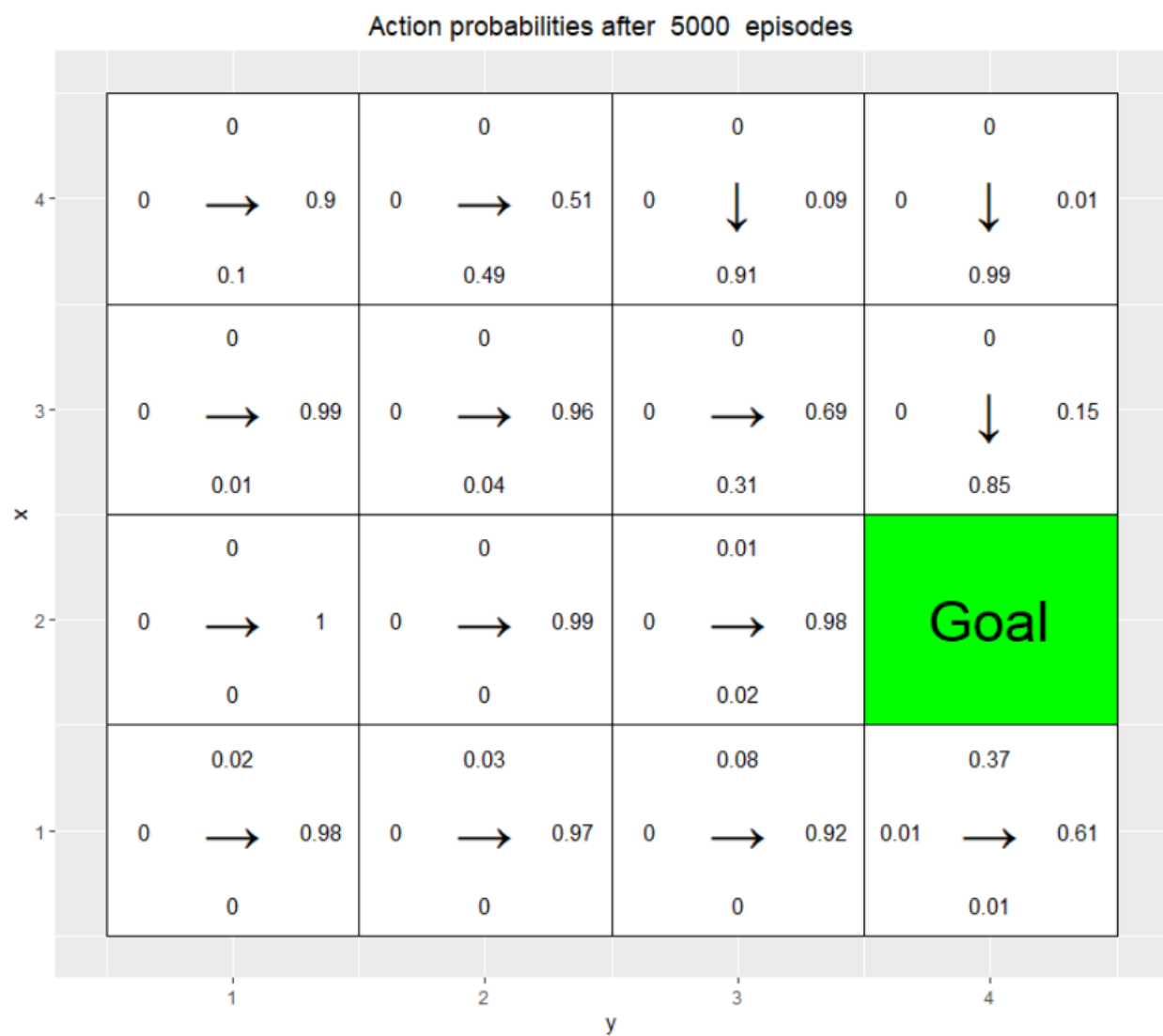


Figure 12: Reinforce with goal in (2,4)



Figure 13: Reinforce with goal in (1, 1)