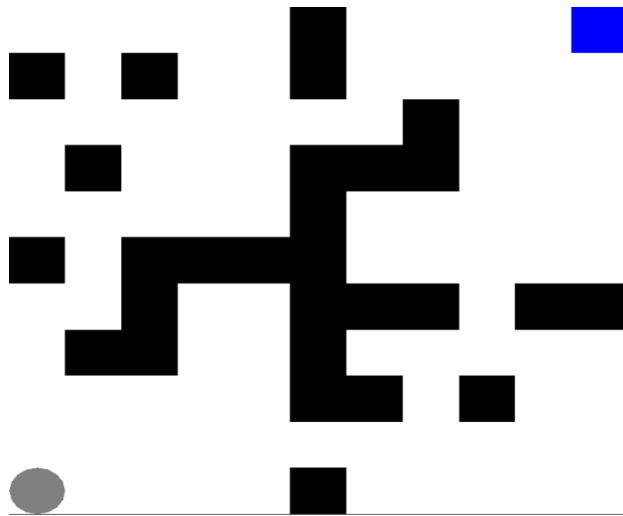


## Assignment 4 Report

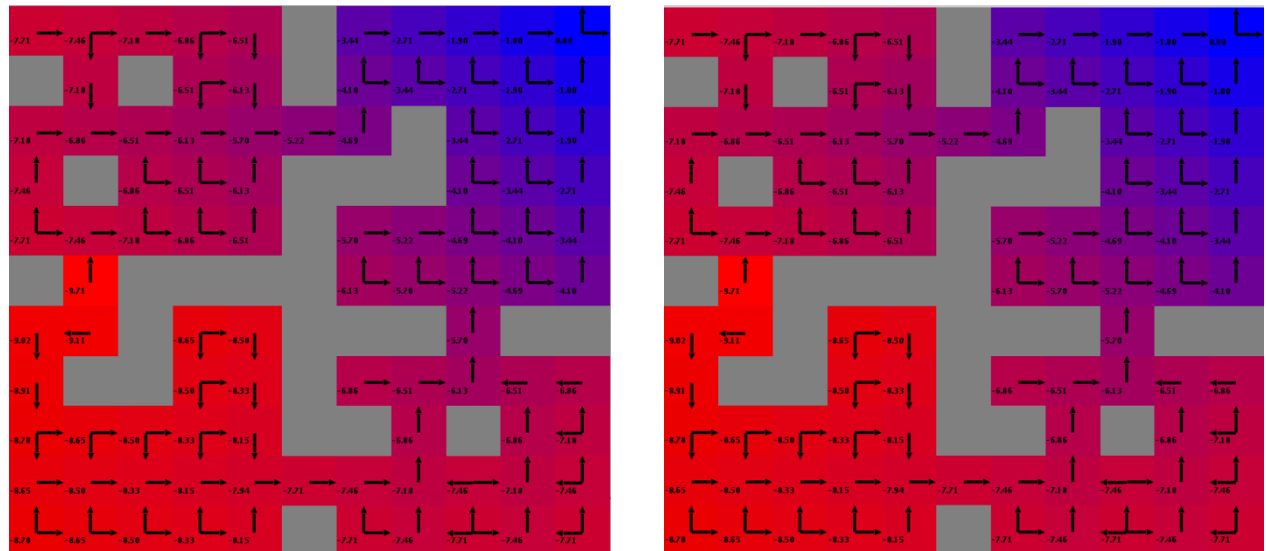
### MDP Description:

Both MDPs are simply gridworlds, where a single agent attempts to find the most optimum policy in traversing through the world. The difference is one has a much larger number of grids, meaning more states to traverse through. This problem can be interesting as a part of training a reinforce learning based agent for videogame AIs. If we can discretize the map of a game we can do reinforcement learning for adaptive enemy or friendly AIs to create a sort of reaction to player actions. Movement would be the first step to this.

### Small Gridworld:

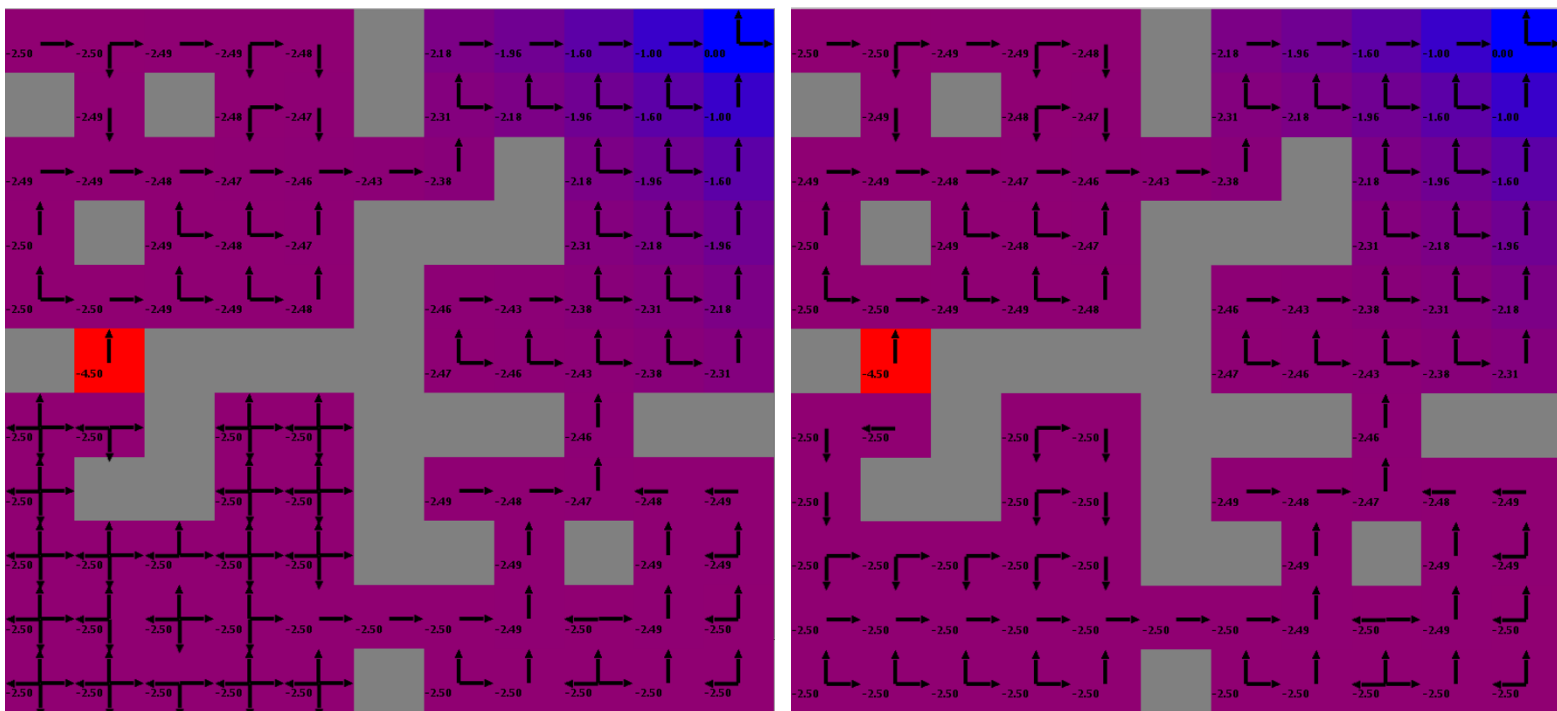


Here is the map of how the gridworld looks. Black tiles are walls, and the blue tile is the terminal state. The rewards given to each tile is -1, and 0 for going out of the terminal state, also the entrance to the left room is given a -3 reward.



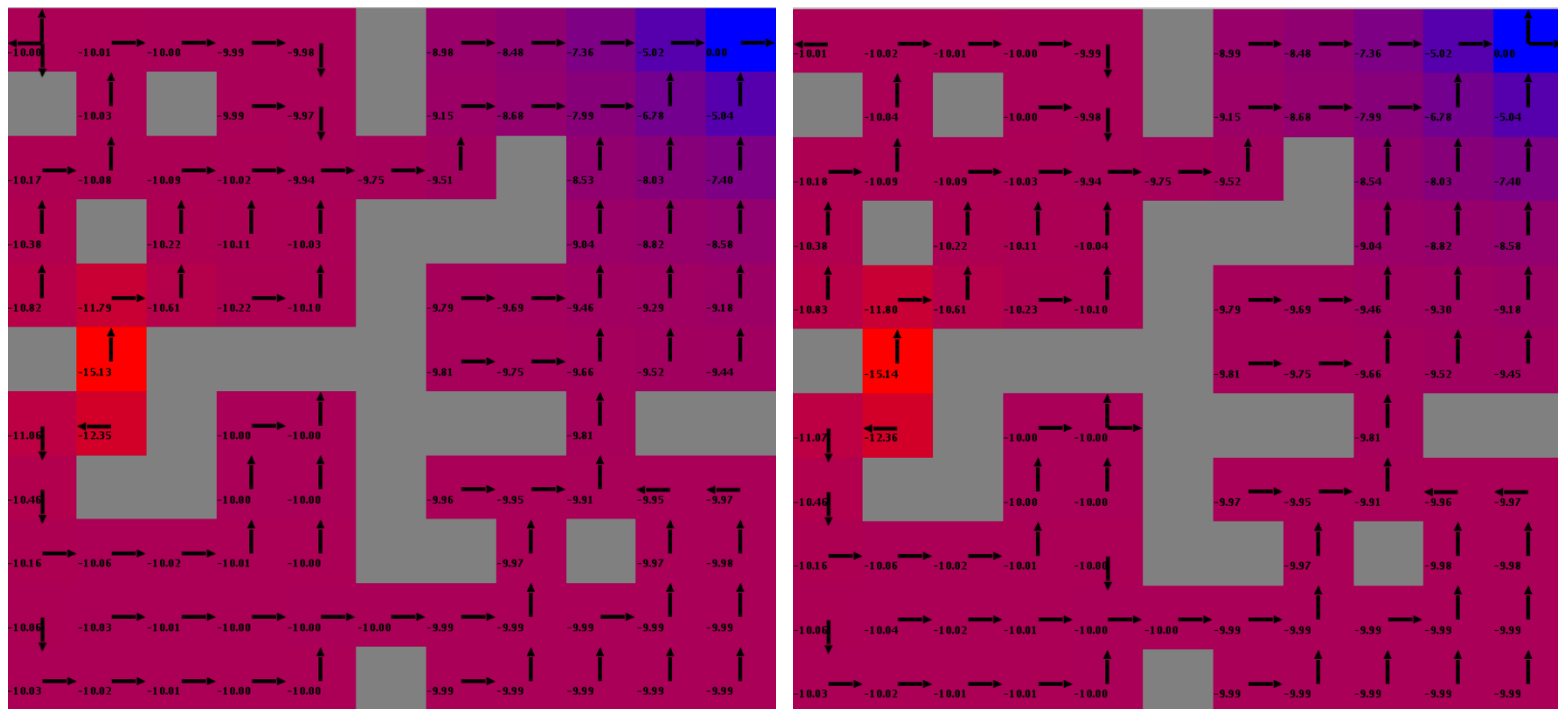
On the left is the resulting policy from value iteration and on the right is from policy iteration. It can be seen that there is no difference in either any of the states values or the policies to take. This means that this is already the optimum policy and the algorithms just terminated themselves as no more improvements can be made. Iteration wise, policy iteration did a lot more iterations with 87 iterations, while value iteration did 23 iterations. This is due to the fact that policy iteration did multiple policy improvement each of which runs many iterations

of policy iteration. Also note that this used a large gamma value, which is 0.9, from my experiments, as I lowered the gamma more and more the policy gets less and less optimal. However with the same gamma values, policy iteration has a more optimum policy compared to value iteration.



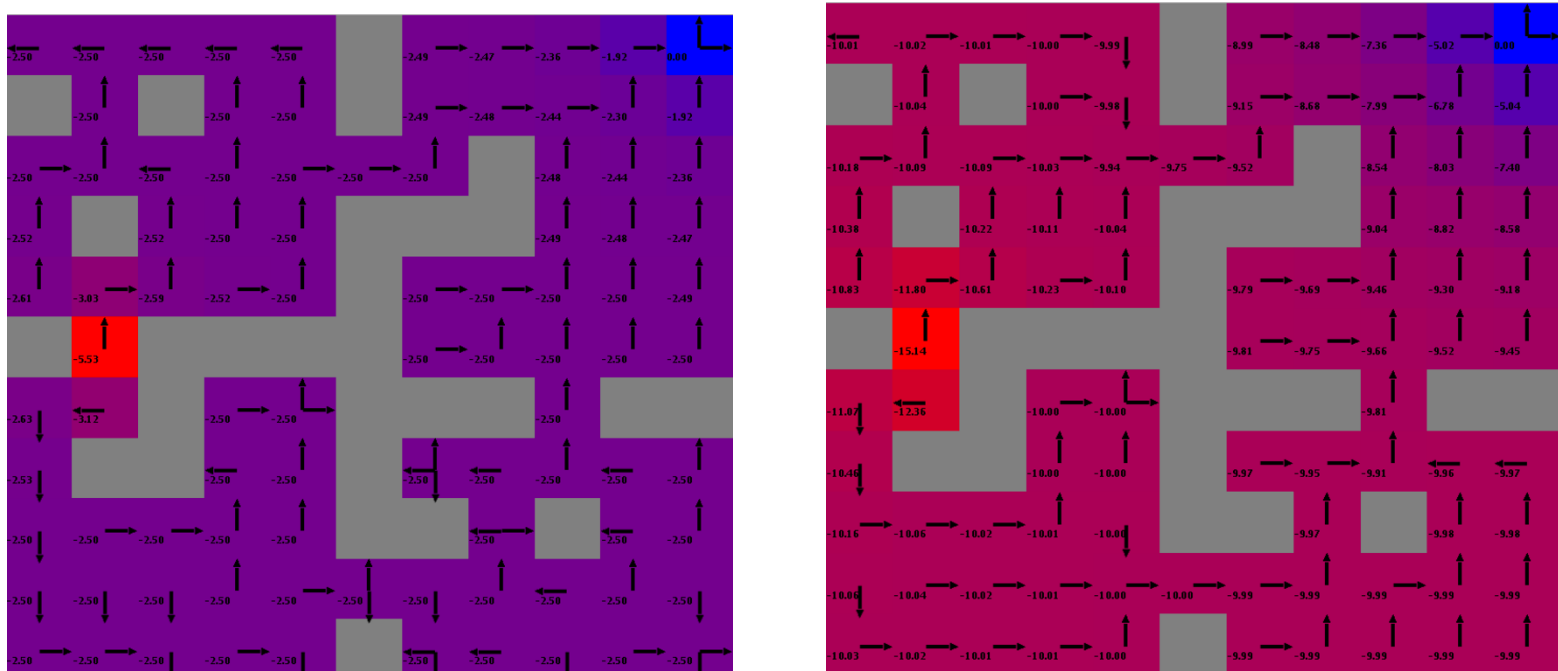
As can be seen, the value iteration results are a lot more confused when it is farther from the reward state. This is due to the reward state being too far and the possible 0 reward being too negligible. Thus near the initial point, the states all seem equally as bad. However due to the policy updating nature of policy iteration, we can still see a clearer policy. Policy iterations still need a larger number of iterations, however for both methods, the amount of iterations needed are far reduced compared to when the gamma is high. This is because the value change is much lower due to the small gamma number and small reward numbers, thus causing it to converge to a value faster. Which can be solved by altering the difference in rewards in the reward function to give a larger living cost (currently living cost is 1, can increase to a larger value). This would result in even larger iteration numbers for large gamma values as the values keep on propagating, resulting in it converging slower.

The pictures below denote a situation where I modified the transition function to not be as certain. The tests above were done when the transition function has a 1.0 probability of moving between tiles. Now when the transition probability is lowered to about 0.3 we can see a difference in some of the policies

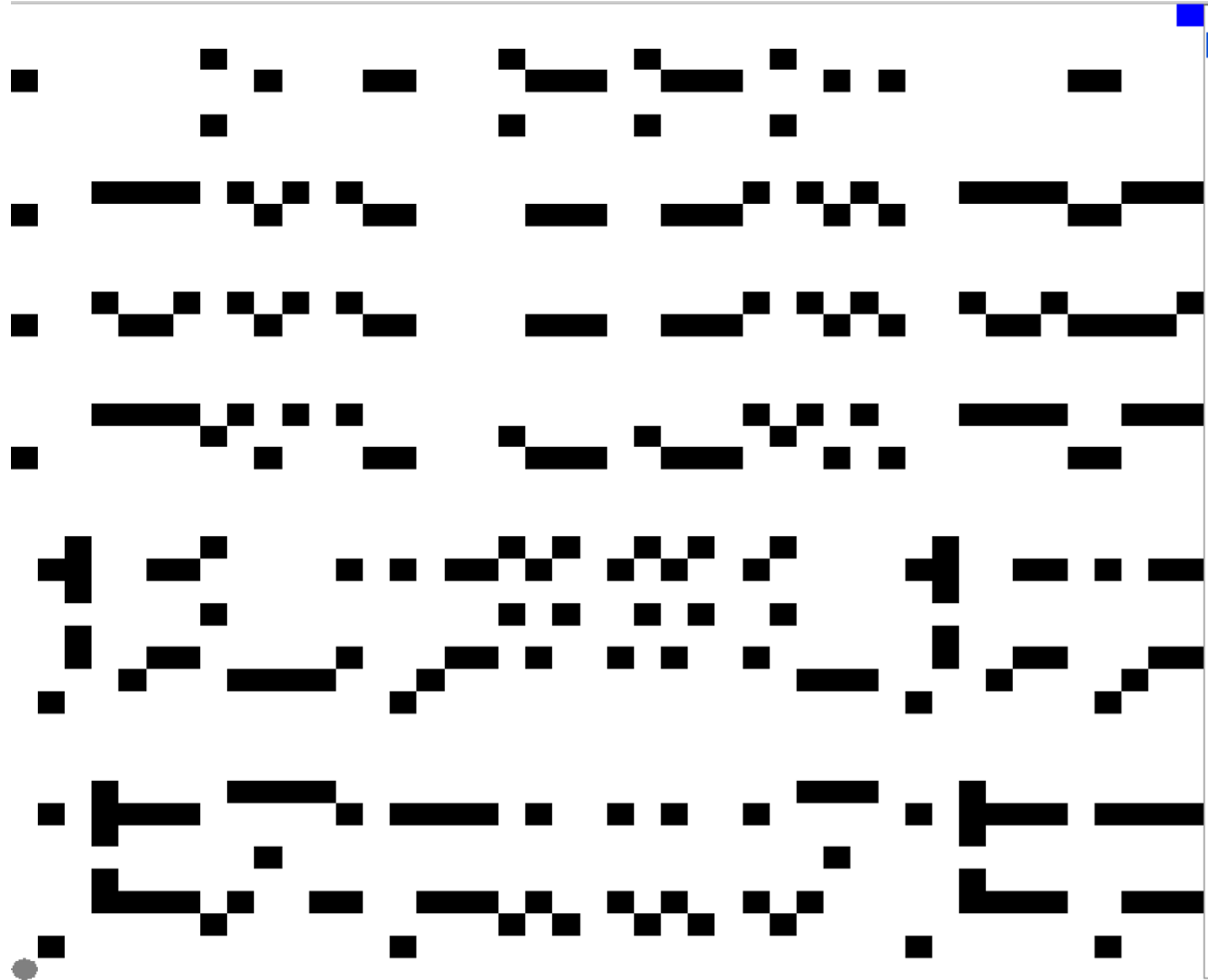


Now here we can still see the policies being fairly similar, but noticeably in states farther from the goal state, the policies are making the agent not go towards the closest exit of the room. It makes the agent move away from the room, this is likely to be caused by the larger probability of the movement failing and the agent ending up in a position where it wasn't going to. However as the terminal state moves closer, the agent simply attempts to reach that terminal state as fast as possible due to the living cost, and is not that affected as the chance of actually ending up where it wants to be in (0.3) is larger than ending up in one of the other states (0.233).

When the discount is adjusted to a smaller number, a similar thing happens with the previous experiment, where the agent will not be able to see how good the terminal state is that is too far away. The behaviors of PI and VI are the same.



## Large Gridworld



In comparison this is the large gridworld's map, where there are a larger number of grids and thus the agent has to go around further and more states for the algorithms to keep track of. Note that this gridworld is actual 16 times larger than the small gridworld.

First things first is that there is a considerable amount of time increase for both policy and value iteration, this took around 10 times longer than the small gridworld. Also another thing is that the reward had to be readjusted to a larger value so that the algorithm can actually converge. With a smaller reward value, most of the states have a really really miniscule difference between states and as such the agent just sees every single state as equally as bad, and this is a convergent policy as the values and policies have already converged to the resulting final values.

-98-98-98.00-98-98-97-97.27-97-96-96.96-95-95-94.24-93-92-92.03-90-89-87-86.83-83-81-79.77-74-71-68.63-61-56-52-46.80-34-27-19.00-0.01  
-99-98-98.00-98-98-98-97.97-97-97-96.96-96-95-95.28-94-93-92-92.03-90-89-87.84-84-83-81.47-77-74-71.50-65-61-56-52.46-40-34-27.00-1.0  
-99-99-98-98.00-98-98.36-97-97-97-97.20-96-96-95.26-94.77-93-92.02-91.14-89.06-86-84.99-81-79-77-74.58-68-65-61-56.53-46-40-34.29-1.9  
-99-99-98.92-98-98-98.36.18-97-97.90.22-96.98-95-95.76-92-92.82-86-84.99.32-79.41-74-71-68-65-61.56.95-40.56-27  
-99-99-99-99.00-98-98-98.92-98-97-97.27-97-96-96.96-95-94-94.89-92-92-91.89-89-87-86.89-83-81-79.47-74-71-68-65.63-56-52-46.80-34  
-99-99-99-99.93-98-98.80-98-98-98-97.97-97-96.96-96.18-95-94.24-93.54-92.03-90-89.06-86-84-83-81.47-77-74-71-68.63-61-56-52.46-40  
-99-99-99-99.20-99-98-98.00-98-98-98.97-97-97-97.20-96-96-95-95.28-94-93-92.03-91-90-89.06-86-84-83.83-79-77-74-71.50-65-61-56.53-46  
-99-99-99-99.20-99-99-98.92-98-98-98.36-97-97-97.20-96-96-96-95.26-94-94-93.94-92-91-90.89-87-86-84-83.83-79-77-74.58-68-65-61.56-52  
-99-99-99.43-99.03-98.80-98.52-98-97-97.27-97-96-96-96.96-95-94-94.89-92.82-91.14-89.06-86-84.99.32-71-68.62  
-99-99-99-93-99-99-99.83.21-98-98.00.80-97-97.27-97.22-95-95.29-92-92.02.82-91.14-86.89-86-87-89.00.15-94.89-95  
-99-99-99-99.80-99-99-99.20-99-98-98.00-98-96-98-97.97-97-96.96-96-95-95.28-94-93-92.03-91-90-89-87.84-87-85-90.89-92-92-93.94-94  
-99-99-99-99.94-99-99-99.20-99-99-98.92-98-96-98-97.97-97-97.20-96-96-95.26-94-94-93.94-92-91-90-89.06-89-90-91.84-92-93-94.89-95  
-99-99-99-99.90-99-99-99.36-99-99-99.00-98-96-98-97.97-97.20-96-96-96.96-95-94-94.89-92-92-91-90.89-90-91-92.02-93-94-94.75-95  
-99-99-99.66-99-99.54-99.43-99.29-99.13-98-96-98-98.92-98-97-97.27-97-96-96.96-96.96-95.76-94.77-93.54-92-91.84.15-92-92.82-94-94-95.26.76  
-99-99-99.73-99-99.90.54-99-99.20.29-98-98.00-98.36-97-97.22-95-95.26.76-94.77-92.02-92.02-95.29-97  
-99-99-99-99.20-99-99-99.94-99-99-99.20-99-99-98-98.00-98-98-98.97-97-97.20-96-96-95.26-94-94-93-92.02-92-93-94.89-95-95-96.96-96  
-99-99-99-99.93-99-99-99.20-99-99-99.36-99-99-99.00-98-96-98-97.97-97.20-96-96-96-95.26-94-94-93.94-93-94-94.75-95-96-96.96-97  
-99-99-99-99.20-99-99-99.00-99-99-99.36-99-99-99.00-98-98-98.92-98-97-97.20-96-96-96-96.96-95-94-94.89-94-94-95.26-96-96-96.96-97  
-99-99-99.80-99.66-99.58-99.48-99-99-99-99.83-98-98-98.00-98-98-97.97-97.50-96.91-96.18-95-94.24.19-96-96.91  
-99-99-99-99.84-99-99.78-99-99-99-99.80-99-99-99.20.13-98-98.00-98.52-98.96-97-97.97-96.96-96-95-95.28-95-95-96.96-96-97-97.20-97  
-99-99-99.02-99-99-99.20.70-99-99-99.54-99-99-99.29-98-98.80-97-97-97.22-96.57-95.26-95-96-96.96.91-97.27-98  
-99-99-99-99.84-99-99-99.20-99-99-99.36-99-99-99.00-99-99-99.00-98-98-98.92-98-97-97-97.20-96-96-96.96-96-96-96.96-97-97.97-97.96-98  
-99-99-99-99.94-99-99-99.20-99-99-99.00-99-99-99.36-99-99-99.00-98-98-98.92-98-98-98.92-97-97-97.20-96-96-96.96-96-97-97.20-97-98-98-98.92  
-99-99.90-99.00-99-99.84-99-99-99-99.20-99-99-99.20.48-99.36-99-99.13-98.92-98-98.52-98.96-97-97-97.22-97-97-97.27-98-98-98-98.00  
-99.92-99.98.88-99-99.80-99-99.75-99.70-99.62-99.58-99.30.29-99.29-98.98.80-98-98.36-97-97.75-97-97.75-98.36-98.67  
-99-99.94-99.90-99-99-99.80-99-99-99.36-99-99-99.00-99-99-99.00-99-99-99.00-98-98-98.92-98-98-98.92-97-97-98.18-97-97-98.92-98-98-98-98.92  
-99-99-99-99.90-99-99.88-99.80-99-99-99.20-99-99-99.00.62-99.54-99-99.36-99.21-99-98.92-98.00-98-98-98.36-97-96-98.36-98-98-98-99.00  
-99-99.94-99.92-99-99-99-99.80-99-99-99.00-99-99-99.20-99-99-99.94-99-99-99.20-99-98-98.00-98-98-98.52-98-98-98.92-98-98-99-99.00  
-99-99.94-99.93.92-99-99.80-99-99.84-99-99-99.75-99.66-99.30.54-99.43-99.29.21-99-98.92-98-98-98.67-98-98.52-98-99-99.13  
-99-99-99-99.94-99-99-99.90-99.84-99.80-99.20-99-99-99.00-99-99-99.36-99-99.13-98-98.00-98.52-99-99.83-99.13-99-99  
-99-99-99-99.94-99-99-99-99.80-99-99-99.00-99-99-99.84-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-98-99.00-98-96-99-99.83-99-99-99.83  
-99-99-99-99-99.94-99-99-99.80-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.20-99-99-99.83-98-96-99-99.83-99-99-99.83  
-99-99-99-99-99.94-99-99-99.80-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.20-99-99-99.83-98-96-99-99.83-99-99-99.83  
-99-99-99-99-99.94-99-99-99.80-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.20-99-99-99.83-98-96-99-99.83-99-99-99.83  
-99-99-99.97-99.90-99-99.94-99.90-99-99-99.80-99-99-99.20-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.20  
-99.98-99.97-99-99.90-99-99.94-99.90-99.84-99-99.78-99.73-99.00.62-99-99.90-99-99.29-99.29-99.48  
-99-99-99.97-99.92-99-99-99.80-99-99-99.00-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.94-99-99-99.83.36-99-99-99.80-99-99-99.00  
-99-99-99-99-99-99-99-99.95-99-99-99.93-99-99-99.93-99-99-99.00-99-99-99.00-99-99-99.20-99-99-99.58-99-99-99-99.83-99-99-99.00-99-99-99.20  
-99-99-99.98-99-99-99-99-99.90-99-99-99.94-99-99-99.90-99-99-99.84-99-99-99.20-99-99-99.00-99-99-99.48-99-99-99.00-99-99-99.20  
-99-99-99.98-99.97-99-99.95-99-99.93-99.89-99-99.85-99.82-99.20.75-99.70-99-99-99-99.54-99-99.70  
-99.99-99-99-99.98-99.90-99-99-99.93-99-99-99.93-99.91-99.00.88-99.85-99-99.80-99.20-99-99-99.90-99-99-99.20-99-99-99.20  
-99.99-99-99-99.98-99-99-99.90-99-99-99.93.95-99-99-99.90-99-99-99.00-99-99-99.00-99-99-99.20-99-99.62-99.02-99-99-99.20.73-99-99.00  
-99.99-99-99-99.98-99-99-99.90-99-99-99.90-99-99-99.93-99-99-99.00-99-99-99.84-99-99-99.20-99-99-99.00-99-99-99.20-99-99-99.00

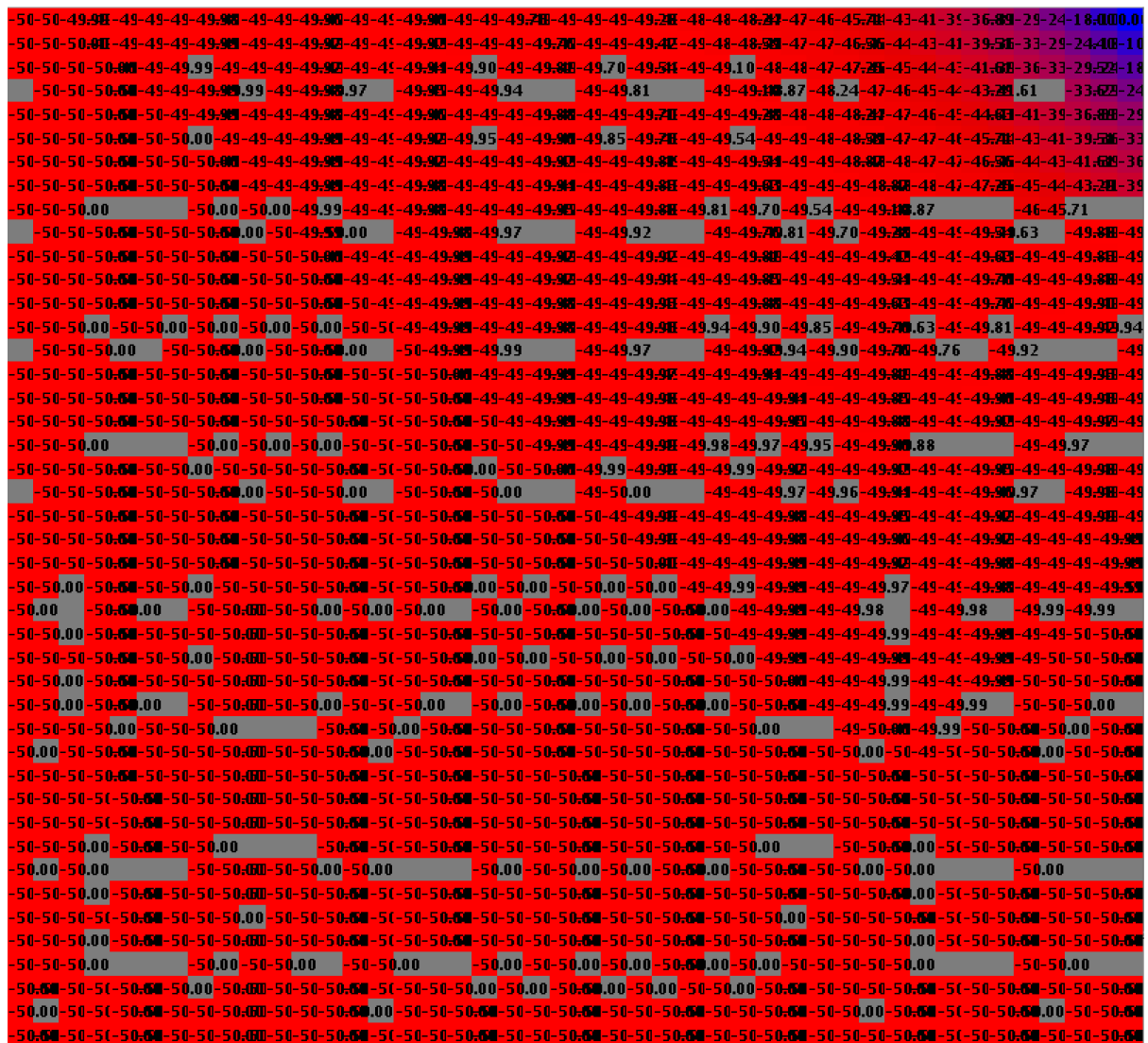
Here there is a considerably much larger amount of states such that I am not able to fully show the policies for every single state but on the above is the results of the value iteration and the values denoted by it, here the agent will simply attempt to do a beeline straight from its' starting position to the terminal state as there are no major obstacles along the way. Which is expected given the reward function which is just a living cost, some small rewards scattered along the grid (which is just a smaller living cost), and the 0 terminal state.



[illegible]

[illegible]





And lastly, the above are the respective graphs for VI and PI when the discount factor is lowered (from 0.9 to a smaller number), notably even with a small decrease, from 0.9 to 0.8, the reward amounts scale down really quickly, resulting in reaching convergence faster due to the really small difference. Thus again causing the issue of each state being equally as bad as its surroundings except when close to the terminal state.

Conclusion:

So a larger state space causes problems in the sense that the reward function needs to be varied as the difference between states will get smaller and smaller as the number of states increase. This can result in less than optimal results due to the implementation of the training algorithms. This can usually be solved by increasing the living cost (or reward of the terminal state depending on implementation, this gridworld is implemented such that when the agent enters the terminal state, it always gets 0 reward, thus to increase the reward of entering the terminal state, we must lower the reward of every other state). The discount factor also matters in the case of how good the policies are, a too low discount will cause the agent to focus too much on local rewards, and a larger discount will help it focus on a more global reward system.