

CS 7641 - Markov Decision Processes

Author: Swapnil Ralhan(sralhan3@gatech.edu)

Date: Nov 30, 2014

Abstract

In this paper, we explore Markov Decision Processes. We take 2 MDPs, one with a small(~10) number of states and a large(~600) number of states. To solve these problems, we compare 3 approaches - Value Iteration, Policy Iteration and Q-Learning. We consider how long each of the approaches take to converge, the solutions they converge to, and different approaches used within each problem(different reward functions and terminal functions, as well different values for the various variables).

Problems

The following are two variations of the traditional grid problem. They are interesting as they easily highlight the different features and applicability of the RL and MDP algorithms, as well as being simple to understand and easy to visualize. Additionally, they also find direct applications in path-finding problems, such as robotics, and AI in computer games.

Small Problem

Size

The first problem we look at is a 3 by 3 grid

Obstacles

The grid has 2 obstacles - one at (0,2) and another at (1,0).

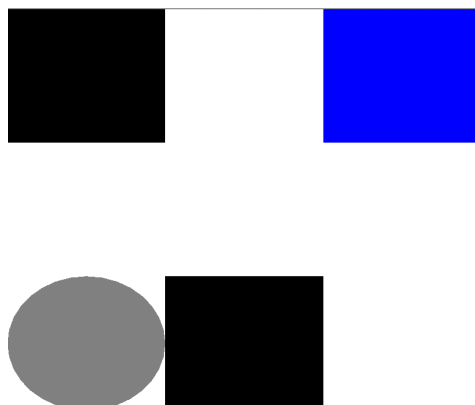
Starting Positions

The Agent starts at (0,0) and the goal is at (2,2).

Reward Function

Each step was uniform cost. Entering the goal state resulting in termination, thus the agent wanted to end the game as soon as possible. The cost of each step was 1.0

The representation of the grid looks as follows:



Large Problem

Size

The second problem we examine is a 24 by 24 grid.

Obstacles

There are 6 columns of obstacles - At x index 2, 10 and 18 we have a column from y index 0 to 22, and at y indices 6, 14 and 22, we have a column from y index 1 to 23.

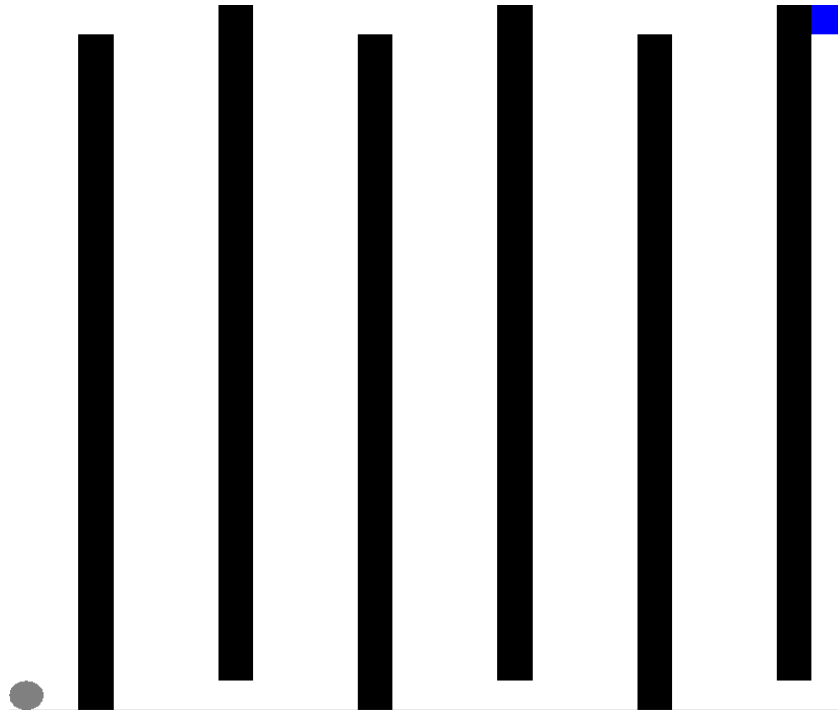
Starting Positions

The Agent starts at (0,0) and the goal is at (23, 23).

Reward Function

Each step was uniform cost. Entering the goal state resulting in termination, thus the agent wanted to end the game as soon as possible. The cost of each step was 1.0

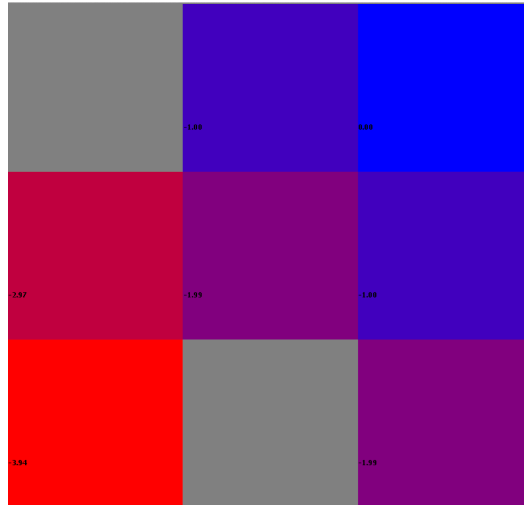
The grid looks as follows:



Problem 1 - Small Grid

Value Iteration

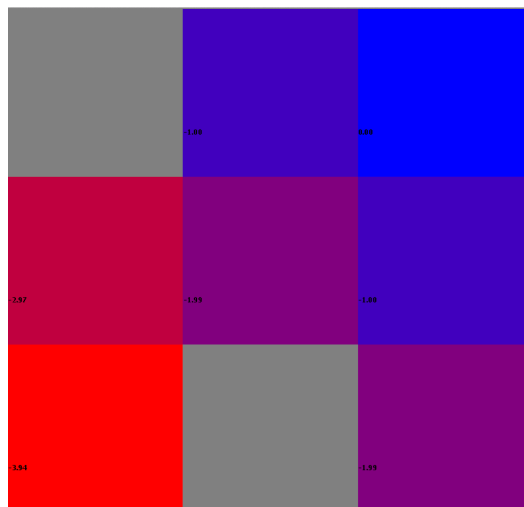
Value Iteration solved the problem within 4 Passes. The values assigned to each of the states are shown in the diagram below.



The solution obtained was optimal, taking 4 steps to the solution state.

Policy Iteration

Policy Iteration also solved the problem in 4 iterations, with the values assigned to each state being the same as those assigned by Value Iteration:

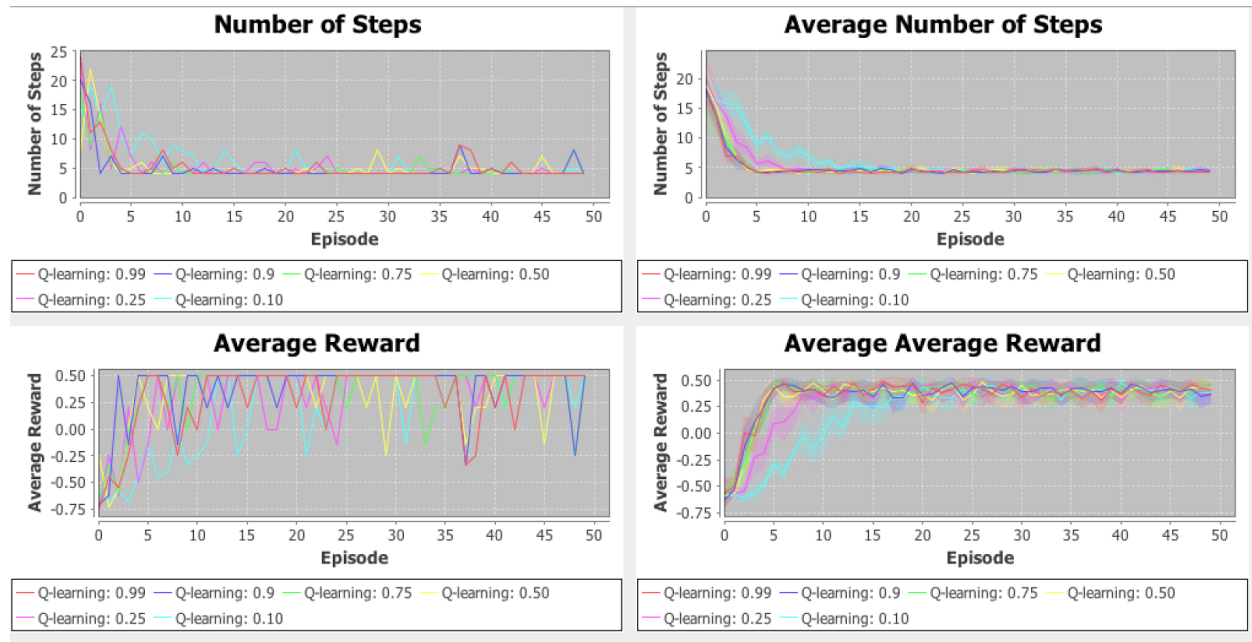


The solution obtained was also optimal.

Q-Learning

Q-learning converged in approximately 2 iterations. The solution obtained is optimal as well (taking 4 steps).

Trying out different values for the discount rate, we see that lower discount values slow down the learning, however, the optimal is achieved within 20 iterations for every single value of gamma.



Comparison

Among the 3 algorithms, each of them have very similar performance, reaching convergence with a few iterations. While Q-learning is marginally faster in terms of the number of iterations, VI iterations run much faster than either Q-learning or Policy Iteration.

Problem 2 - Large Grid

Value Iteration

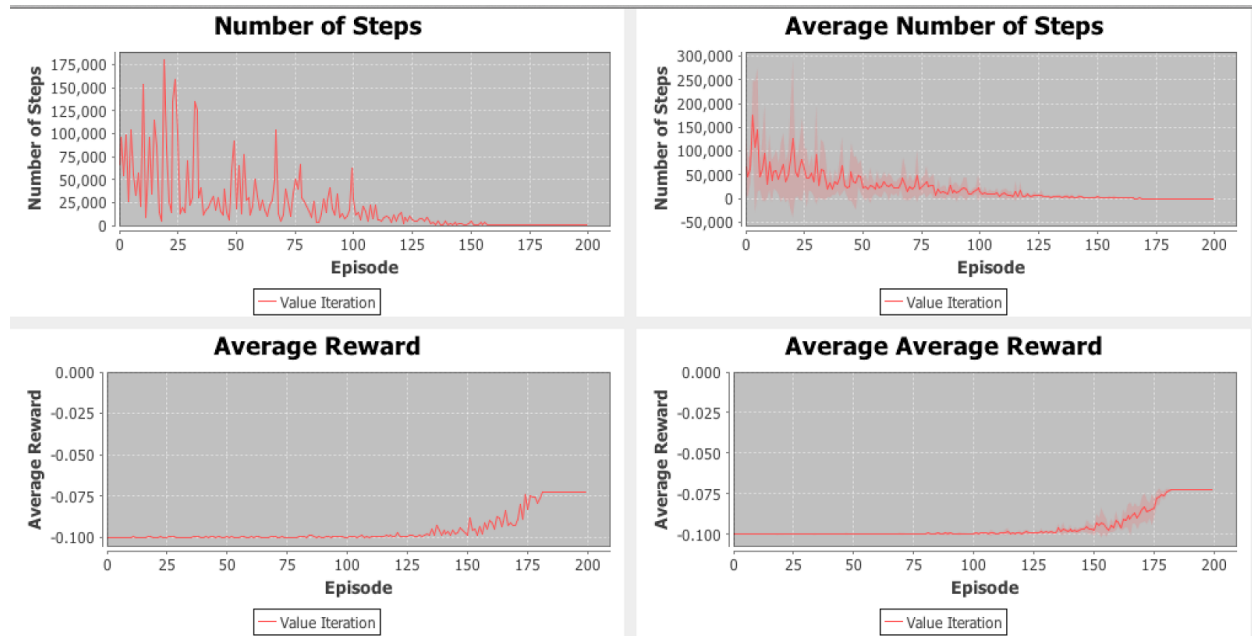
Value Iteration takes 184 Iterations to find the optimal solution of length 184. The values produced for each location are shown below:

-80.17~-79.97~-79.77~-79.57~-79.37~-79.15	-66.27~-65.87~-65.57~-65.17~-64.87~-64.47~-64.13	-41.87~-41.37~-40.77~-40.17~-39.57~-38.87~-38.27	0.00			
-80.37~-80.17	-79.37~-79.17~-78.97	-66.57~-66.27~-65.87	-64.47~-64.17~-63.76	-42.47~-41.87~-41.30	-38.87~-38.27~-37.65	-1.00
-80.57~-80.37	-79.17~-78.97~-78.73	-66.97~-66.57~-66.22	-64.17~-63.77~-63.40	-43.07~-42.47~-41.88	-38.27~-37.67~-37.02	-1.99
-80.77~-80.57	-78.97~-78.77~-78.51	-67.27~-66.97~-66.56	-63.77~-63.47~-63.03	-43.67~-43.07~-42.46	-37.67~-37.07~-36.38	-2.97
-80.97~-80.76	-78.77~-78.57~-78.30	-67.57~-67.27~-66.90	-63.47~-63.07~-62.65	-44.17~-43.67~-43.04	-37.07~-36.37~-35.74	-3.94
-81.17~-80.95	-78.57~-78.37~-78.08	-67.87~-67.57~-67.23	-63.07~-62.67~-62.28	-44.77~-44.17~-43.61	-36.37~-35.77~-35.09	-4.90
-81.37~-81.14	-78.37~-78.07~-77.85	-68.27~-67.87~-67.56	-62.67~-62.27~-61.90	-45.27~-44.77~-44.17	-35.77~-35.07~-34.43	-5.85
-81.57~-81.33	-78.07~-77.87~-77.63	-68.57~-68.27~-67.88	-62.27~-61.97~-61.51	-45.87~-45.27~-44.73	-35.07~-34.47~-33.77	-6.79
-81.77~-81.52	-77.87~-77.67~-77.41	-68.87~-68.57~-68.20	-61.97~-61.57~-61.12	-46.37~-45.87~-45.28	-34.47~-33.77~-33.10	-7.73
-81.87~-81.70	-77.67~-77.47~-77.18	-69.17~-68.87~-68.52	-61.57~-61.17~-60.73	-46.97~-46.37~-45.83	-33.77~-33.17~-32.43	-8.65
-82.07~-81.89	-77.47~-77.17~-76.95	-69.47~-69.17~-68.83	-61.17~-60.77~-60.33	-47.47~-46.97~-46.37	-33.17~-32.47~-31.74	-9.56
-82.27~-82.07	-77.17~-76.97~-76.71	-69.77~-69.47~-69.15	-60.77~-60.37~-59.93	-47.97~-47.47~-46.91	-32.47~-31.77~-31.06	-10.47
-82.47~-82.25	-76.97~-76.71~-76.48	-70.07~-69.77~-69.45	-60.37~-59.97~-59.53	-48.47~-47.97~-47.44	-31.77~-31.07~-30.36	-11.36
-82.67~-82.43	-76.77~-76.47~-76.24	-70.37~-70.07~-69.76	-59.97~-59.57~-59.12	-49.07~-48.47~-47.97	-31.07~-30.37~-29.66	-12.25
-82.77~-82.60	-76.47~-76.27~-76.00	-70.67~-70.37~-70.06	-59.57~-59.17~-58.71	-49.57~-49.07~-48.49	-30.37~-29.67~-28.94	-13.13
-82.97~-82.78	-76.27~-76.07~-75.76	-70.97~-70.67~-70.36	-59.17~-58.77~-58.29	-50.07~-49.57~-49.00	-29.67~-28.97~-28.23	-13.95
-83.17~-82.95	-76.07~-75.77~-75.51	-71.27~-70.97~-70.66	-58.77~-58.27~-57.87	-50.57~-50.07~-49.51	-28.97~-28.27~-27.50	-14.85
-83.27~-83.12	-75.77~-75.57~-75.27	-71.57~-71.27~-70.95	-58.27~-57.87~-57.44	-51.07~-50.57~-50.02	-28.27~-27.57~-26.77	-15.71
-83.47~-83.29	-75.57~-75.27~-75.02	-71.87~-71.57~-71.24	-57.87~-57.47~-57.01	-51.57~-51.07~-50.52	-27.57~-26.77~-26.03	-16.55
-83.67~-83.45	-75.27~-75.07~-74.76	-72.17~-71.87~-71.53	-57.47~-57.07~-56.58	-51.97~-51.57~-51.01	-26.77~-26.07~-25.28	-17.38
-83.77~-83.62	-75.07~-74.77~-74.51	-72.37~-72.17~-71.81	-57.07~-56.57~-56.14	-52.47~-51.97~-51.50	-26.07~-25.27~-24.53	-18.21
-83.97~-83.78	-74.77~-74.57~-74.25	-72.67~-72.37~-72.10	-56.57~-56.17~-55.70	-52.97~-52.47~-51.99	-25.27~-24.57~-23.77	-19.03
-84.17~-83.95	-74.57~-74.27~-73.99	-72.97~-72.67~-72.37	-56.17~-55.77~-55.25	-53.47~-52.97~-52.47	-24.57~-23.77~-23.00	-19.87
-84.27~-84.11	-74.27~-73.97~-73.77~-73.47~-73.27~-72.97~-72.65	-55.77~-55.27~-54.87~-54.37~-53.87~-53.41~-52.94	-23.77~-23.07~-22.27~-21.47~-20.67			

These values are produced by having a cost of -1 per step, and the terminal position being the only way to end.

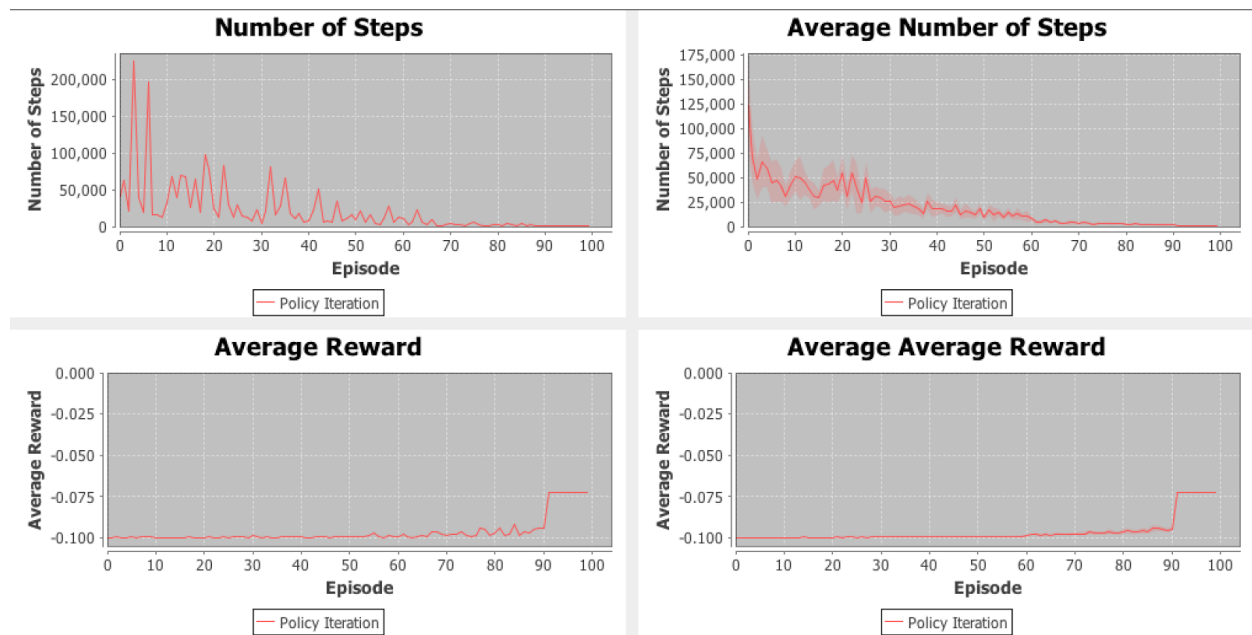
In the following graphs, we show the number of steps suggested by value iteration to get to the solution state, as a factor of the number of iterations taken by VI. As you can see, it starts off by randomly exploring the space, taking up 250,000 steps in some cases. Eventually, as the number of iterations increases, the solution begins to converge. At exactly 184 iterations, it has found the optimal solution.

In terms of the average reward, as we can see, the reward starts rising steeply at around 160 iterations, peeking at the optimal value obtained at 184 iterations.



Policy Iteration

Policy Iteration takes much fewer iterations. Within 93 Iterations, Policy Iteration is able to find the optimal solution of 184 steps. However it should be noted here that the 93 iterations takes about twice as long as the 184 iterations taken by Value Iteration(5.26 seconds as opposed to 2.50 seconds).



-80.17~-79.97~-79.77~-79.57~-79.37~-79.15	-66.27~-65.87~-65.57~-65.17~-64.87~-64.47~-64.13	-41.87~-41.37~-40.77~-40.17~-39.57~-38.87~-38.27	0.00			
-80.37~-80.17	-79.37~-79.17~-78.94	-66.57~-66.27~-65.88	-64.47~-64.17~-63.76	-42.47~-41.87~-41.30	-38.87~-38.27~-37.65	-1.00
-80.57~-80.37	-79.17~-78.97~-78.73	-66.97~-66.57~-66.22	-64.17~-63.77~-63.40	-43.07~-42.47~-41.88	-38.27~-37.67~-37.02	-1.99
-80.77~-80.57	-78.97~-78.77~-78.51	-67.27~-66.97~-66.56	-63.77~-63.47~-63.03	-43.67~-43.07~-42.46	-37.67~-37.07~-36.38	-2.97
-80.97~-80.76	-78.77~-78.57~-78.30	-67.57~-67.27~-66.90	-63.47~-63.07~-62.65	-44.17~-43.67~-43.04	-37.07~-36.37~-35.74	-3.94
-81.17~-80.95	-78.57~-78.37~-78.08	-67.87~-67.57~-67.23	-63.07~-62.67~-62.28	-44.77~-44.17~-43.61	-36.37~-35.77~-35.09	-4.90
-81.37~-81.14	-78.37~-78.07~-77.85	-68.27~-67.87~-67.56	-62.67~-62.27~-61.90	-45.27~-44.77~-44.17	-35.77~-35.07~-34.43	-5.85
-81.57~-81.33	-78.07~-77.87~-77.63	-68.57~-68.27~-67.88	-62.27~-61.97~-61.51	-45.87~-45.27~-44.73	-35.07~-34.47~-33.77	-6.79
-81.77~-81.52	-77.87~-77.67~-77.41	-68.87~-68.57~-68.20	-61.97~-61.57~-61.12	-46.37~-45.87~-45.28	-34.47~-33.77~-33.10	-7.73
-81.87~-81.70	-77.67~-77.47~-77.18	-69.17~-68.87~-68.52	-61.57~-61.17~-60.73	-46.97~-46.37~-45.83	-33.77~-33.17~-32.43	-8.65
-82.07~-81.89	-77.47~-77.17~-76.95	-69.47~-69.17~-68.83	-61.17~-60.77~-60.33	-47.47~-46.97~-46.37	-33.17~-32.47~-31.74	-9.56
-82.27~-82.07	-77.17~-76.97~-76.71	-69.77~-69.47~-69.15	-60.77~-60.37~-59.93	-47.97~-47.47~-46.91	-32.47~-31.77~-31.06	-10.47
-82.47~-82.25	-76.97~-76.77~-76.48	-70.07~-69.77~-69.45	-60.37~-59.97~-59.53	-48.47~-47.97~-47.44	-31.77~-31.07~-30.36	-11.36
-82.67~-82.43	-76.77~-76.47~-76.24	-70.37~-70.07~-69.76	-59.97~-59.57~-59.12	-49.07~-48.47~-47.97	-31.07~-30.37~-29.66	-12.25
-82.77~-82.60	-76.47~-76.27~-76.00	-70.67~-70.37~-70.06	-59.57~-59.17~-58.71	-49.57~-49.07~-48.49	-30.37~-29.67~-28.94	-13.11
-82.97~-82.78	-76.27~-76.07~-75.76	-70.97~-70.67~-70.36	-59.17~-58.77~-58.29	-50.07~-49.57~-49.00	-29.67~-28.97~-28.23	-13.95
-83.17~-82.95	-76.07~-75.77~-75.51	-71.27~-70.97~-70.66	-58.77~-58.27~-57.87	-50.57~-50.07~-49.51	-28.97~-28.27~-27.50	-14.85
-83.27~-83.12	-75.77~-75.57~-75.27	-71.57~-71.27~-70.95	-58.27~-57.87~-57.44	-51.07~-50.57~-50.02	-28.27~-27.57~-26.77	-15.71
-83.47~-83.29	-75.57~-75.27~-75.02	-71.87~-71.57~-71.24	-57.87~-57.47~-57.01	-51.57~-51.07~-50.52	-27.57~-26.77~-26.03	-16.55
-83.67~-83.45	-75.27~-75.07~-74.76	-72.17~-71.87~-71.53	-57.47~-57.07~-56.58	-51.97~-51.57~-51.01	-26.77~-26.07~-25.28	-17.34
-83.77~-83.62	-75.07~-74.77~-74.51	-72.37~-72.17~-71.81	-57.07~-56.57~-56.14	-52.47~-51.97~-51.50	-26.07~-25.27~-24.53	-18.21
-83.97~-83.78	-74.77~-74.57~-74.25	-72.67~-72.37~-72.10	-56.57~-56.17~-55.70	-52.97~-52.47~-51.99	-25.27~-24.57~-23.77	-19.03
-84.17~-83.95	-74.57~-74.27~-73.99	-72.97~-72.67~-72.37	-56.17~-55.77~-55.25	-53.47~-52.97~-52.47	-24.57~-23.77~-23.00	-19.81
-84.27~-84.11	-74.27~-73.97~-73.77~-73.47~-73.27~-72.97~-72.65	-55.77~-55.27~-54.87~-54.37~-53.87~-53.47~-52.94	-23.77~-23.07~-22.27~-21.47~-20.67			

Again, we see that the values for each of the states produced by Policy Iteration are exactly the same as those produced by Value Iteration.



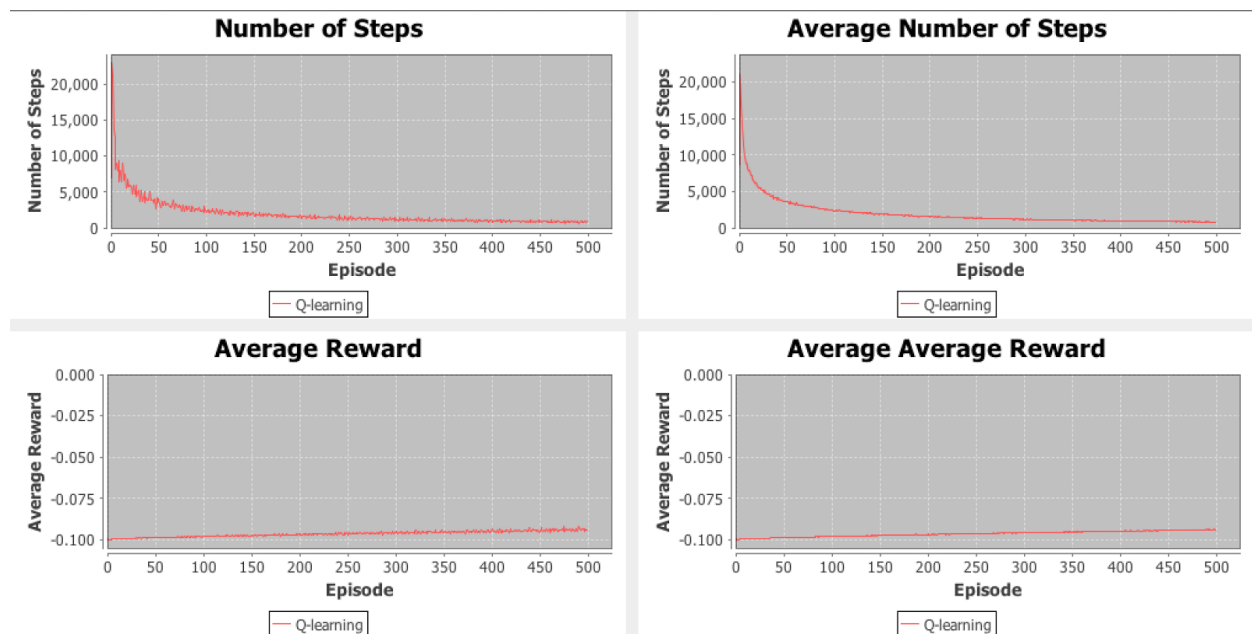
Here we investigated different values of gamma(discount factor). At values of 0.99, 0.95, the same optimal values are obtained, albeit at a slightly higher number of iterations. However, for discounts of 0.90, 0.80 and 0.50, we are very far from the optimal value(in the 1000s of steps), and we are unable to get to the optimal value within 100 iterations.

Q-Learning

Now, we tried Q-Learning on the larger problem. Trying out with a discount rate and a learning rate of 0.9, we see convergence in approximately 600 iterations. However, the iterations take a shorter time than Value Iteration or Policy Iteration, taking about 4.54 seconds to do the 600 iterations.

The convergence achieved is at 184 steps, the optimal value.

Here we plot the number of steps as a function of the episodes:



Next we experiment with different values of gamma, the discount rate. As we can see, for lower discount rates, q-learning sometimes fails to converge. Specifically, we can see that for gamma=0.50, we start diverging at about 250 episodes. The minimum value achieved is on the order of 10000 steps for gamma=0.50.

For higher values of gamma, we eventually converge, but taking more iterations for smaller gammas:



Next we investigate different learning rates. As we can see, different learning rates cause convergence to take longer. A learning rate of 0.25 causes convergence in over 800 episodes, and a learning rate of 0.10 does not come close to convergence in over 1000 episodes.

It should be noted though that even though we see reward for learning rate of 0.9 top off at around 300 episodes, this is at a suboptimal path of approximately 190-210 steps. Since the value of -1 is assigned to each step, the reward difference between 200 steps and 184 steps is minimal



We also explored different reward functions. We tried assigning values between 0 and 5000 to the goal state, and between -1 and -0.1 per step. Placing particular emphasis on the ratio of step cost vs goal reward, we did not see any significant change in the performance of the q-learning algorithm.

Conclusion

Thus we have evaluated the performance of the 3 algorithms - Value Iteration, Policy Iteration and Q-learning on the 2 different problems. As we have seen, on both the problems each of the algorithms achieves the optimal step value. While Policy Iteration takes a fewer number of iterations for the larger problem, Q-learning converges in a smaller number of iterations on the smaller problem. However, iteration time is different as well, and Value Iteration trumps the other algorithms when it comes to time taken to achieve the optimal value.

We have also seen, that choosing the right number of minimum iterations, learning rate and discount rate(γ) is very important as well. We may fail to converge(if choosing < 100 iterations for Value Iteration), or even diverge(if choosing too low a discount rate for Q-learning for instance). Thus, it would be an interesting exercise to plot these values for much larger number of episodes, to see if we see any divergence in the future. For instance, if we stopped Q-learning at $\gamma=0.50$ at 200 episodes, we would mistakenly conclude that the algorithm has converged.

An additional follow-up could be to investigate different problems of the gridworld, using a fixed(such as the alternating walls of the large problem). We can see the evolution of the tuning parameters as well as the number of iterations for convergence.