

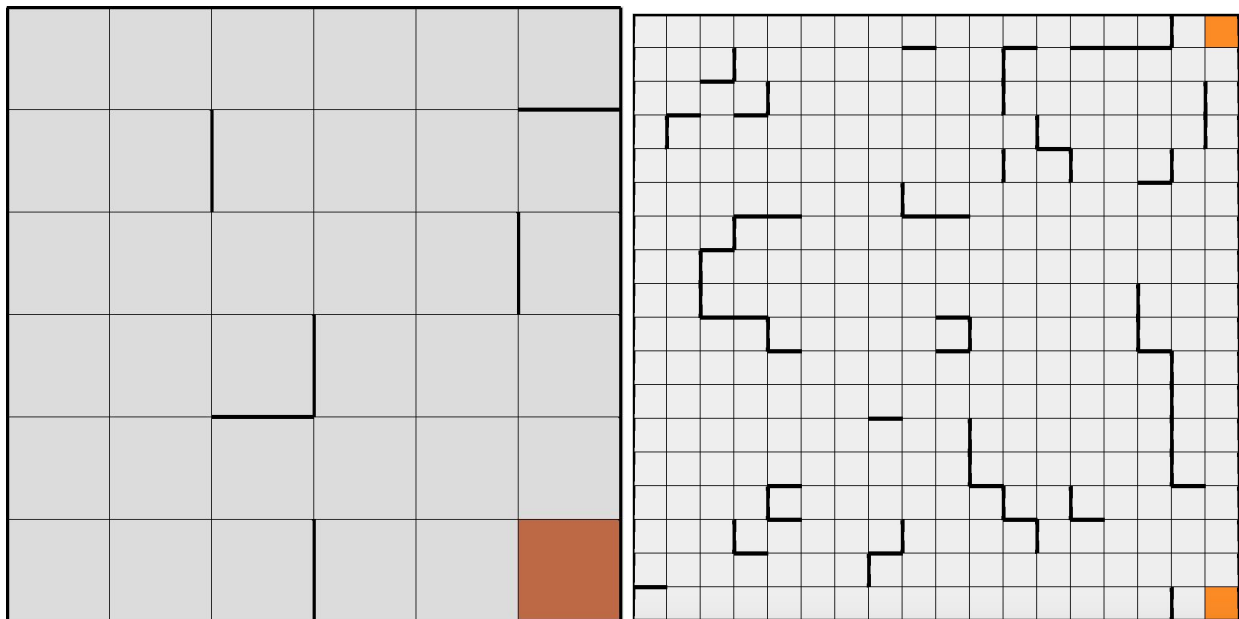
Jeffrey Minowa

CS 4641

Markov Decision Processes

## Problems Used

There were two MDP problems that were created to display the differences in performance of the algorithms. Walls were created by having a -50 value when crossing them, and goals were created by giving the algorithm a reward when found. The simple 6x6 grid had fewer walls, 36 different states, and one goal state. This problem was fairly small and a contrast to the larger 18x18 grid with 324 states, many more walls, and multiple goal states. These two are good to juxtapose because they display different samples for the algorithms to iterate. With the smaller grid, both policy and value iteration were quick to converge while avoiding walls. With the larger grid, it was interesting to see how the algorithm made preferences toward the multiple goal states, given the placement of the walls.



The hyperparameter tested was PJOG which represented the percent chance of going a suboptimal path. This introduces the idea of exploration vs exploitation. Exploration's goal is to find other paths to take in order to find potential optimal, unexplored paths. The goal of exploitation is to use paths already found to hone in on the optimal state by using previous information.

## Background on Value Iteration and Policy Iteration

Two similar algorithms that are introduced in this paper are value iteration and policy iteration. Both of the algorithms need to know the transition states prior to actually being run. Each have their pros and cons and can be used to solve Markov Decision Processes (MDPs) which will be discussed later in the paper.

For now, we will define the overarching idea behind each algorithm:

Value iteration, as the name suggests, iteratively recalculates the values of the states until there is convergence.

Policy iteration, on the other hand, computes its policy by solving a set of linear equations from expected reward and iteratively checks to see if other policies would improve performance until a policy is guaranteed to be optimal <sup>[1]</sup>.

## Small Map Comparison

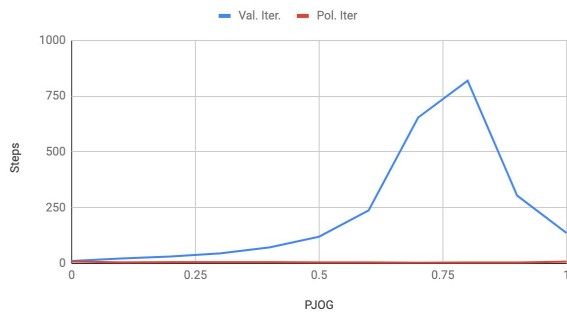
As expected, the number of steps required for the policy iteration is much lower than the value iteration, especially closer to larger PJOG values. As can be seen by the line graphs below, policy iteration consistently has a smaller number of steps needed to converge onto its answer. This is because policy iteration is more efficient in terms of number of iterations needed.

However, policy iteration also has the downside of being computationally more expensive per iteration. This difference is caused by the fact that value iteration continually recalculates the value of being in each state until convergence, whereas, policy iteration recalculates based on policy and not value which results in it having a two step process - calculating its linear equations and finding out if a policy needs changing. Clearly, finding the same policy will be faster than recalculation until a specific confidence (in all cases in the paper, the precision value used was .0001) which is why the number of iterations is consistently low.

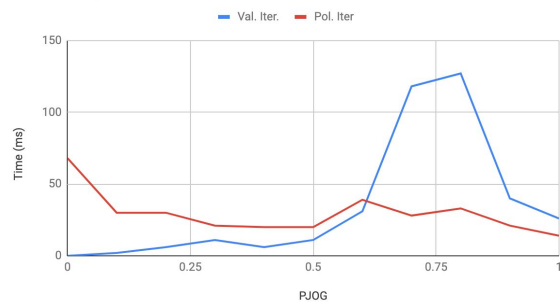
Being a more intensive process in each step also provides insight to why policy iteration takes longer for lower PJOG values. For most of the PJOG values, policy iteration's time is significantly larger until the PJOG is large. This divergence is likely caused by value iteration's bias toward calculating to convergence of the smallest precision. When manually stepping through the algorithm, value iteration changes only one policy every 10s of iterations once it has come close to convergence. Those changing policies are usually in the states that can go either direction. Because the PJOG peaks at .8, that may be the area that value iteration has the closest values between two or more policies, and this would cause values to oscillate instead of converge.

Another noteworthy aspect is that value iteration performs poorly in larger PJOG values because exploring the optimal map and converging is difficult when there is a high chance of randomness. This may be why policy iteration outperforms value iteration on both charts for large PJOG values.

Small Map: Val. Iter. Steps Vs Pol. Iter. Steps

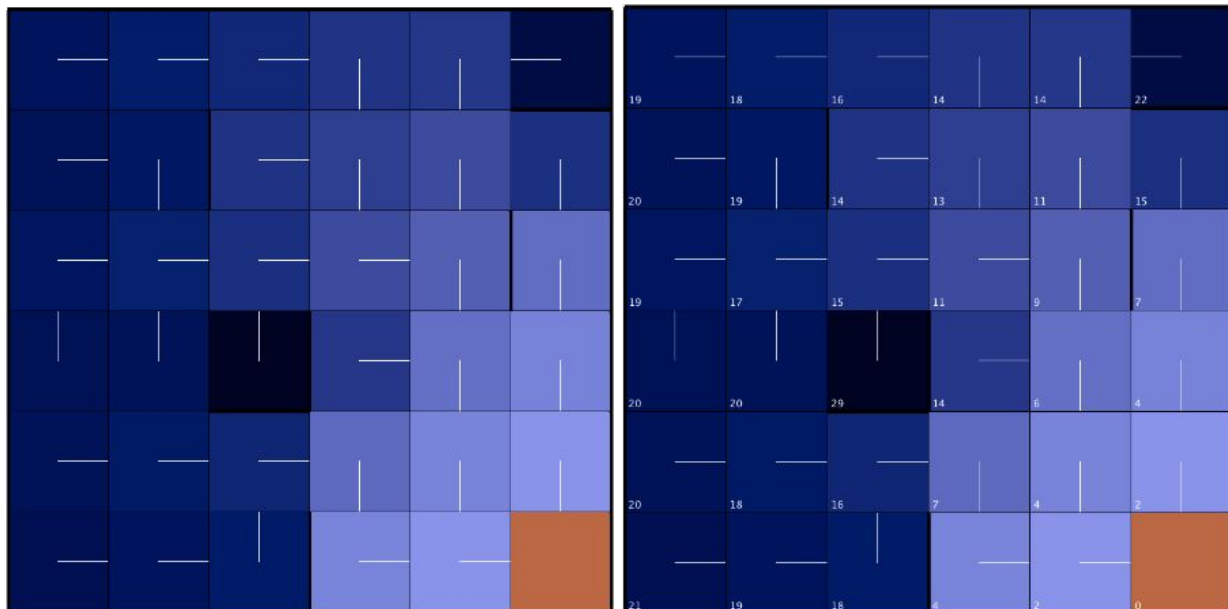


Small Map: Val. Iter. Vs Pol. Iter. Time

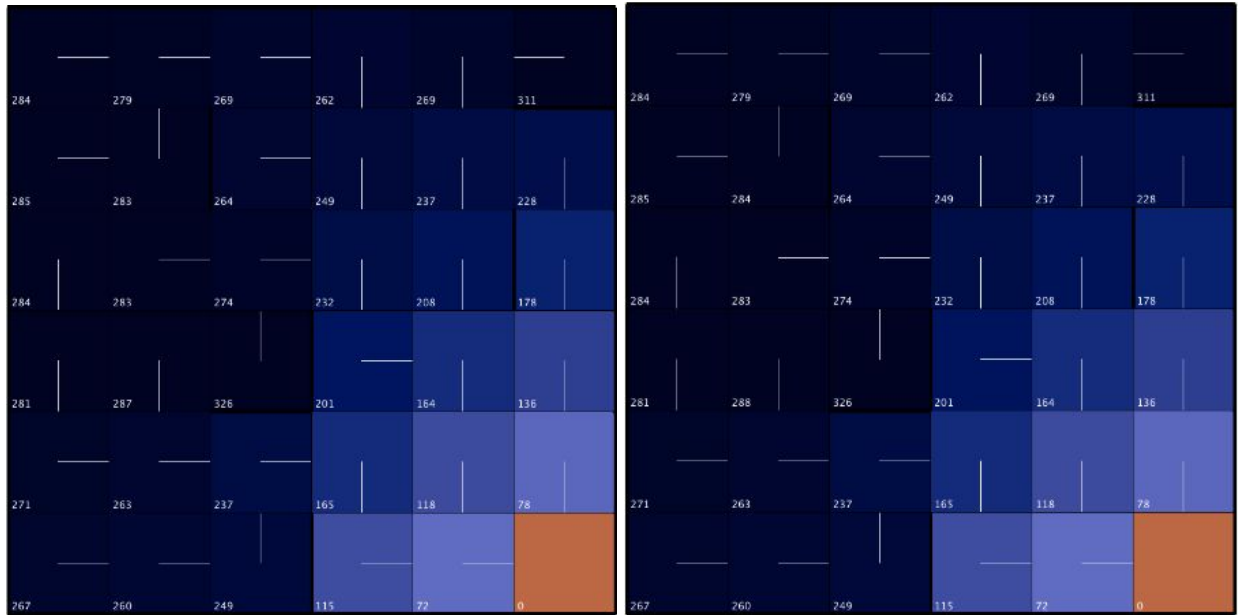


However, for low PJOG values (.3 in the figures below), the policies that both algorithms converge to are exactly the same in the case of the small grids regardless of algorithm.

(Left: Value Iteration, Right: Policy Iteration with PJOG = .3)



(Left: Value Iteration, Right: Policy Iteration with PJOG = .7)

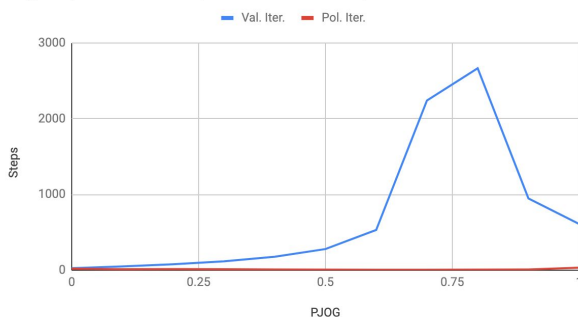


## Big Map Comparison

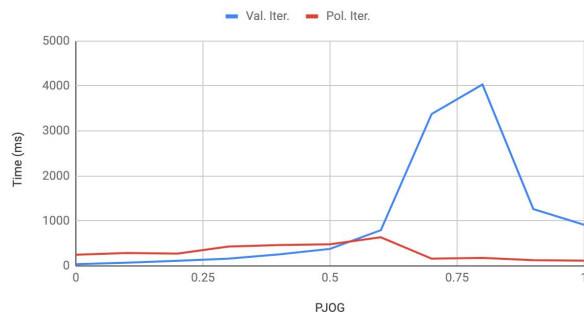
Compared to the small map graphs, the big map's charts for PJOG vs Iterations and PJOG vs Time are similar in shape to the small map's charts. Again, the number of steps for value iteration is magnitudes larger than policy iteration, regardless of the PJOG values. And again, the time the algorithms take to converge cross a little past the .5 mark, displaying that increase in randomness causes value iteration to suffer in performance.

Interestingly, both algorithms have nearly identical policies after convergence even with the two goal states, proving that both find optimal policies to get to the goal states and have a bias toward one goal state over the other. Though it is worth noting that for extremely large grids, policy iteration would take a marginally larger time because as the number of steps increases, so will the time that policy iteration needs to calculate each step.

Big Map: Val. Iter. Steps Vs Pol. Iter. Steps



Big Map: Val. Iter. Steps Vs Pol. Iter. Time



(Left: Value Iteration, Right: Policy Iteration)

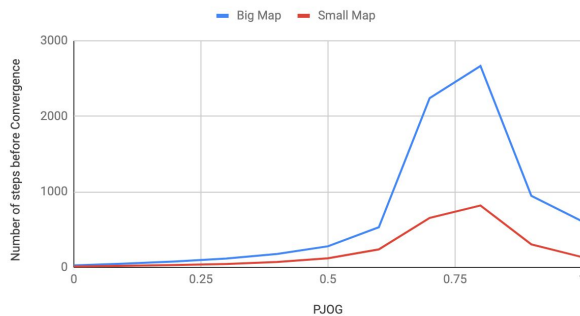
134	120	111	99	90	83	83	84	97	84	84	76	54	77	96	119	16	0
121	112	125	92	80	78	74	74	74	73	80	60	43	43	35	26	16	12
115	106	94	106	80	73	71	71	70	68	74	51	40	35	31	28	26	30
111	97	83	81	73	70	68	69	67	65	63	60	52	38	34	33	36	47
95	84	79	75	71	68	66	68	68	65	63	60	63	45	38	47	46	54
95	86	83	83	76	67	64	68	79	65	54	50	47	44	41	49	49	57
99	90	97	84	73	64	61	60	63	59	51	48	46	45	44	46	49	58
103	92	88	72	65	62	59	57	55	52	49	46	43	43	46	49	52	60
90	84	91	76	66	61	59	58	57	56	48	45	43	43	51	59	56	64
87	78	80	84	76	61	58	56	57	79	50	43	41	40	47	73	61	69
82	74	71	64	68	60	57	54	52	44	41	34	34	34	34	54	70	72
80	71	68	65	63	60	58	58	49	46	42	39	37	33	33	44	70	64
80	71	68	66	64	62	60	60	53	55	47	38	35	33	34	40	63	54
82	73	70	69	69	61	59	56	54	57	55	40	35	32	28	32	58	42
85	76	74	79	86	64	59	56	52	50	58	55	44	40	24	22	24	30
89	80	82	90	76	66	63	71	53	46	44	50	38	31	22	19	17	20
99	85	85	84	74	70	72	65	49	43	39	35	30	25	21	18	13	10
119	96	94	91	83	79	76	67	58	52	48	43	38	34	30	33	15	0

134	120	111	99	90	83	83	84	97	84	84	76	54	77	96	119	16	0
121	112	125	92	80	78	74	74	74	73	80	60	43	43	35	26	16	12
115	106	94	106	80	73	71	71	70	68	74	51	40	35	31	28	26	30
111	97	83	81	74	70	68	69	67	65	63	60	52	38	34	33	36	47
95	84	79	75	71	68	66	68	68	65	63	60	63	45	38	47	46	54
95	86	83	83	76	67	64	68	79	65	54	50	47	44	41	49	49	57
99	90	97	86	73	64	61	60	63	59	51	48	46	45	44	46	49	58
103	92	88	72	65	62	59	57	55	52	49	46	43	43	46	49	52	60
90	84	91	76	66	61	59	58	57	56	48	45	43	43	51	59	56	64
87	78	80	84	76	61	58	56	57	79	50	43	41	40	47	73	61	69
82	74	71	64	68	60	57	54	52	44	41	34	34	34	34	54	70	72
80	71	68	65	63	60	58	58	49	46	42	39	37	33	33	44	70	64
80	71	68	66	64	62	60	60	53	55	47	38	35	33	34	40	63	54
82	73	70	69	69	61	59	56	54	57	55	40	35	32	28	32	58	42
85	76	74	79	86	64	59	56	52	50	58	55	44	40	24	22	24	30
89	80	82	90	76	66	63	71	53	46	44	50	38	31	22	19	17	20
99	85	85	84	74	70	72	65	49	43	39	35	30	25	21	18	13	10
119	96	94	91	83	79	76	67	58	52	48	43	38	34	30	33	15	0

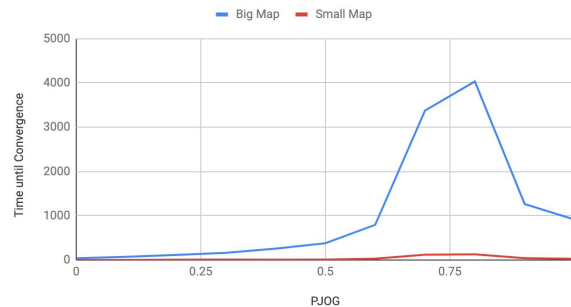
Small vs Big  
 Value Iteration

In regards to having a bigger grid for value iteration, the graphs provide a visual explanation of what changes happened. All the peaks in the small graph are exaggerated by the larger grid. Overall, the algorithm acts as expected when given more states; the algorithm does more work to compensate.

Value Iteration (PJOG vs Steps)



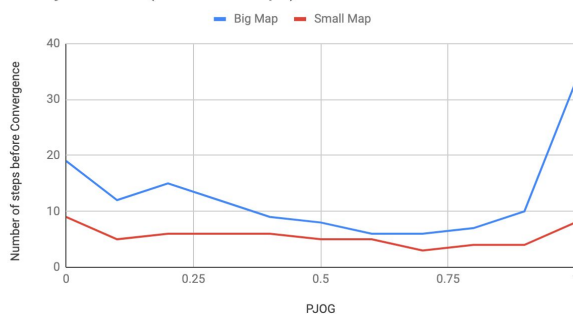
Value Iteration (PJOG vs. Time)



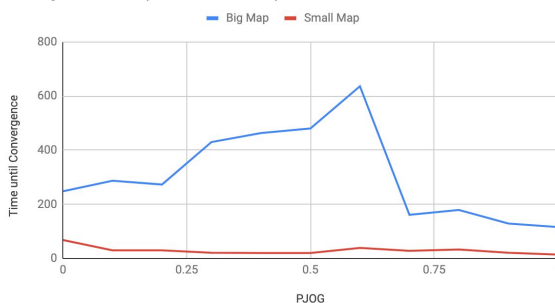
## Policy Iteration

Policy iteration acts differently than value iteration in that it doesn't just scale the graph. A possible explanation on the PJOG vs Time graph divergence is that there is a more even split around .5 and .6, forcing the algorithm to take longer to compute its policies because it is unsure whether to go to the top goal or bottom goal. A possible explanation for the PJOG vs Steps graph's divergence around .9 and 1 is that there may be specific policy that oscillates between directions.

Policy Iteration (PJOG vs. Steps)



Policy Iteration (PJOG vs. Time)



## Q-Learning

Q-learning is another angle to attack MDPs, but in this case, the transition function is unknown. Q-learning is quite different from the previous algorithms because it only knows information that the algorithm has visited and uses that information to make decisions. This particular Q-Learning algorithm also makes use of a decaying learning rate which is similar to simulated annealing's decreasing temperature. The decaying learning rate enforces the bias that early decisions are important in exploration, then goes toward the goal states closest to it. This is important because, unlike the other two algorithms, Q-Learning does not know its transition states, and computes

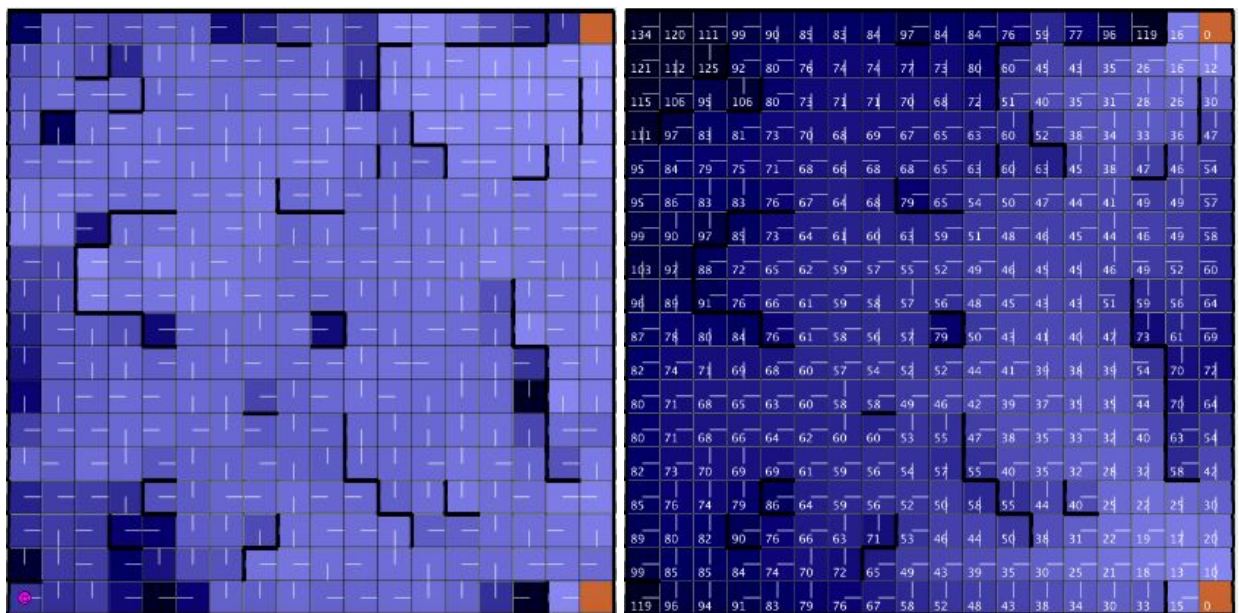


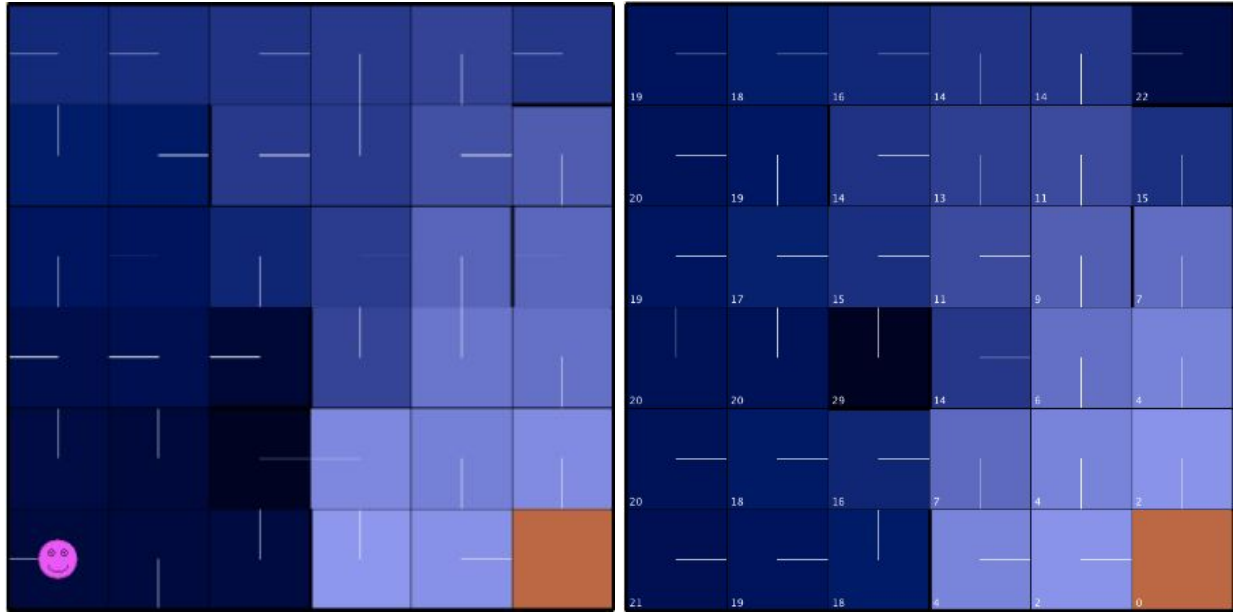
them by exploring different paths. Thus, there is a significant difference in performance between Q-Learning and both Policy Iteration and Value Iteration. Assuming the value iteration has the optimal path, Q-Learning's performance is subpar, and it takes much longer and many more iterations to get to the goal state to learn policies that resemble the optimal. Though this is expected because it is learning as it explores.

The algorithm readily exploits its paths as can be seen from the top left image; Q-learning takes preference toward the top right goal state over the bottom right because of the early stages where the walls are heavily to the right of the starting point, effectively pushing the learner up instead of right.

Big Map			Small Map	
PJOG	0.3		PJOG	0.3
Epsilon	0.1		Epsilon	0.1
Precision	0.001		Precision	0.001
Learning Rate	0.7		Learning Rate	0.7

(Left side: Q Learning, Right side: Value Iteration)





## Conclusion

Because the graphs created from the big map and small map are so similar, these examples consistently reflect the benefits and drawbacks of both algorithms. Value iteration is faster, given exploitation is relatively small. Though policy iteration will take less steps, it is computationally more expensive. Policy iteration is also more robust to randomness as it does not need to converge via strict values. Q-learning is a completely different way to solve the MDP when there is less information given.

### 1. Policy Iteration

<https://www.cs.cmu.edu/afs/cs/project/jair/pub/volume4/kaelbling96a-html/node20.html>