# Project 4: Markov Decision Processes
## CS7641

Anita Rao

arao338@gatech.edu

## 1 MDP PROBLEMS

The two MDP (Markov Decision Process) Problems that we examined in this assignment were (1) Frozen Lake (a grid-world problem) and (2) Forest Management. A goal for the experiments is to study these two MDP problems and compare/contrast the findings between a "small" and "large" problem. We classify the Frozen Lake problem as "small," given that it has very few states (less than 300 max), whereas Forest Management is a "large" problem with 1000+ states.

The Forest Management problem simulates a forest in which forest fires occur with a probability $p$ each year. The **actions** are to 'Wait' (do nothing) or 'Cut' (chop wood). The **reward** is computed based on 1) Cutting wood which generates money (1 pt) and 2) whether we 'Wait' (4 pts) or 'Cut' (2 pts) in the oldest state of the forest. We will explore this problem with **1000 states** (comparison against more states is done below). This problem is interesting due to its large state space and very specific reward function. It may give us insight into how a model-based learner may have an advantage over a model-free learner.

The Frozen Lake problem focuses on a 4x4 grid in which an agent aims to traverse a frozen lake to reach a goal tile -- some squares are walkable and some will result in falling through a hole. The agent will move in the chosen direction with some probability. The **actions** are {Left, Down, Right, Up}. The **reward** is computed based on finding a walkable path to the goal. A lake of size=8x8 and 16x16 was briefly studied as well to understand the behavior of more states. The starting grid for the 4x4 is shown in Figure 1, with 'S' as the start state, 'F' as frozen/walkable, 'H' as a hole, and 'G' as the goal. This compact problem with a tractable reward function may give us insight into whether all 3 algorithms converge on the same optimal solution as expected.



*Figure 1*—Starting state for FrozenLake-4x4 Grid Problem

## 2 ALGORITHMS

The model-based algorithms used to solve the MDP problems above are Value Iteration (VI) and Policy Iteration (PI), where the agent knows the model or reward/transition matrix beforehand and can therefore utilize "offline planning" prior to receiving real-world input. VI uses the Bellman equation to find the optimal value function, and extracts the optimal policy accordingly. PI randomly selects a policy, and continues to improve the policy (by computing the reward function) until the policy no longer changes. QLearning is a model-free algorithm used to solve MDP problems, where the agent *does not* know the model or reward/transition matrix, and therefore does "online" learning as it interacts with its environment. The QLearner for these problems was trained over 1000 iterations for each episode, and the convergence was defined as when the policy stops changing beyond a minimum number of episodes. Another goal of the experiments is to compare/contrast these algorithms for MDP problems.

## 3 EXPERIMENT 1:  VALUE ITERATION

VI was used to evaluate the Forest problem. Figure 2 shows the hypertuning of the **discount factor**, which determines how much the agent factors in future rewards. A smaller discount factor indicates the agent is more greedy, prioritizing immediate rewards. DiscountFactor=0.98 yielded the highest max utility at 44 in 2a, suggesting a greater emphasis on future rewards in determining the optimal value function (**max utility** was calculated by taking the max of the value function at a given iteration). This means prioritizing 'Wait' actions even if it does not yield immediate reward like 'Cut' action. VI converges at 49 iterations for DiscountFactor=0.98 indicated by 2b (**error** was computed by taking the max delta between the current and previous value matrix and termination was reached for error below the library's default value of 0.01).
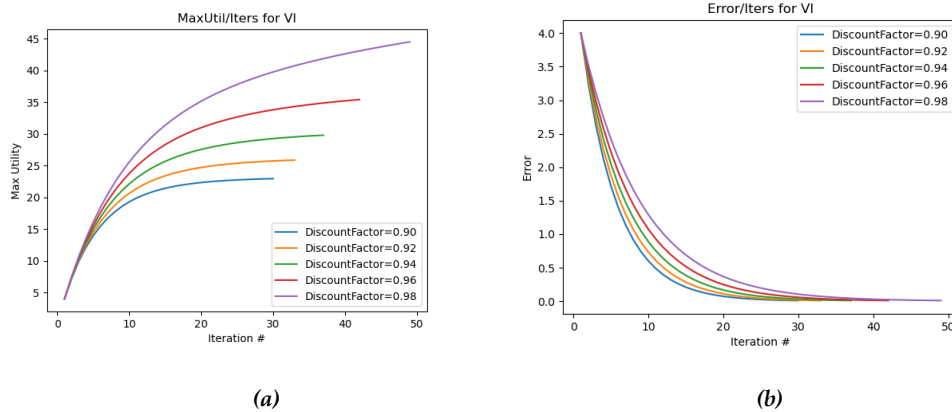


*Figure 2*—*(a)* Forest, VI: Max Utility over Iterations for Discount Factor
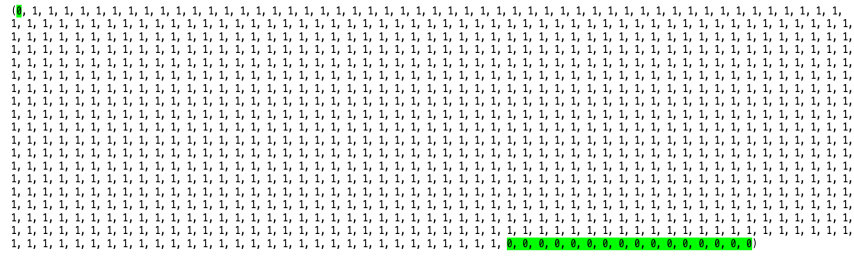*(b)* Forest, VI: Error over Iterations for Discount Factor

*Figure 3*—Forest, VI: Visualization of Policy in Tuple Format

The optimal policy is represented by Figure 3 in tuple format of length 1000 (due to there being 1000 states for the Forest problem), given that this is not a grid-world problem. The policy suggests a 'Wait' in State 0 (after a forest fire has occurred), then a series of 'Cut' to maximize profit, and finally a series of 'Wait' approaching the final, oldest state. This aligns well with what we know about the predefined reward function for the Forest problem.
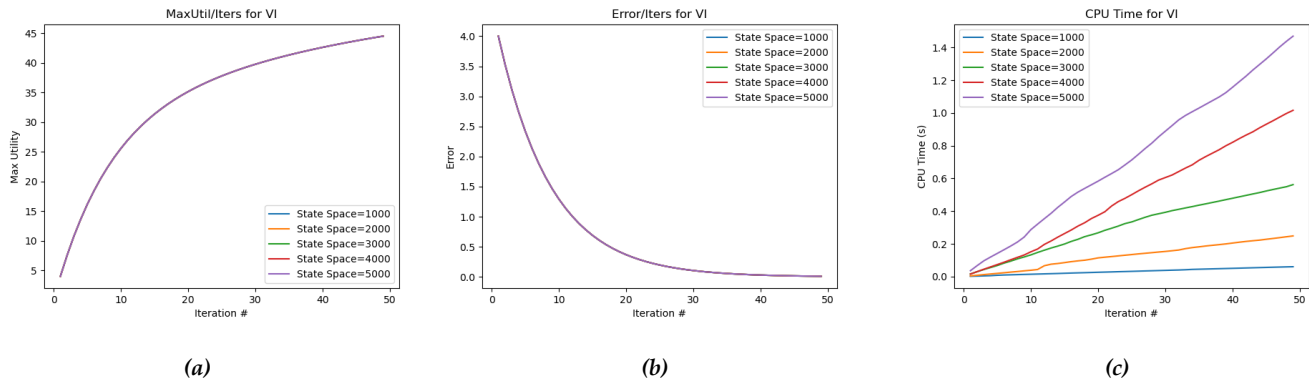


*(a)*                  *(b)*                  *(c)*

*Figure 4*—*(a)* Forest, VI: Max Utility over Iterations for State Space
*(b)* Forest, VI: Error over Iterations for State Space
*(c)* Forest, VI: CPU Time over Iterations for State Space

The state space was also hypertuned to see if more states for the Forest problem had any effect. Figure 4 shows that larger state space did not affect max utility or error in 4a/4b, but naturally increased CPU time in 4c. This provided sufficient confidence that a state space=1000 was adequately big for the Forest Management problem.

Next, VI was used to evaluate the Frozen Lake 4x4 problem. Similar to the Forest problem, a DiscountFactor=0.98 yielded the highest max utility at 0.85 shown in Figure 5a, with a strong emphasis on factoring in future rewards, in this case reaching the goal. VI converged in 33 iterations for DiscountFactor=0.98 based on 5b. The optimal policy is represented in 4c below for FrozenLake-4x4.
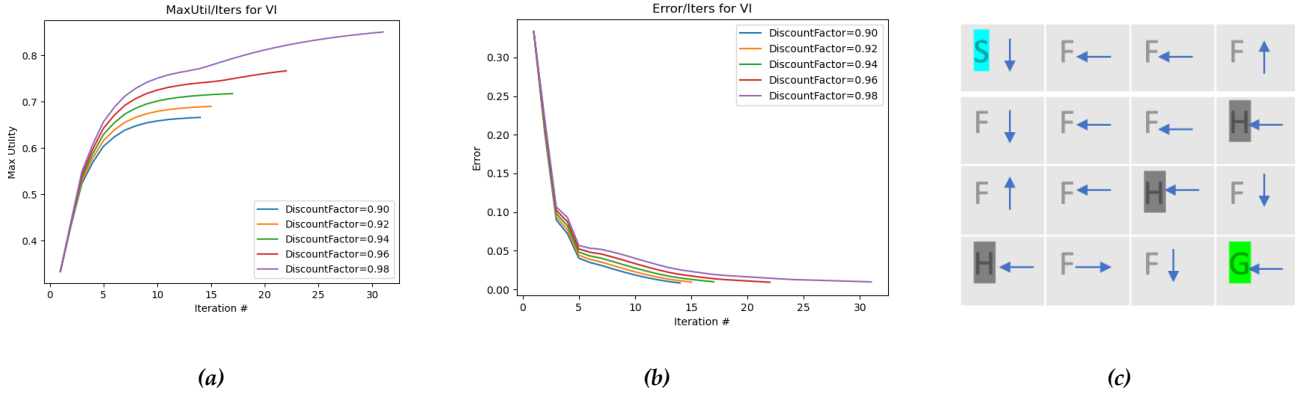
*(a)*     *(b)*     *(c)*

*Figure 5—(a)* FrozenLake, VI: Max Utility over Iterations for Discount Factor
*(b)* FrozenLake, VI: Error over Iterations for Discount Factor
*(c)* FrozenLake, VI: Visualization of Policy in Tuple Format

The grid representation of the 4x4 problem in 5c highlights the optimal policy from VI, overlaid on the original grid in Figure 1. We can see that the suggested pathways are all directed *away* from holes 'H' in the grid. This is a good indication that the agent recognizes no reward associated with the holes. The FrozenLake problem mentions the agent moves in the chosen direction with some probability, so even if the arrow is not pointing towards 'G' the agent can still maneuver to the final goal.

Sweeping the state space for the FrozenLake problem seemed to have varied results, as opposed to what we saw for the Forest problem. 4x4 (16 states), 8x8 (64 states), and 16x6 (256 states) grid problems were compared in Figure 6. There is no obvious trend between the number of states to number of iterations/max utility/compute time. This may be because the Frozen Lake problems we have chosen are very small, causing the curve discrepancies without an underlying trend.
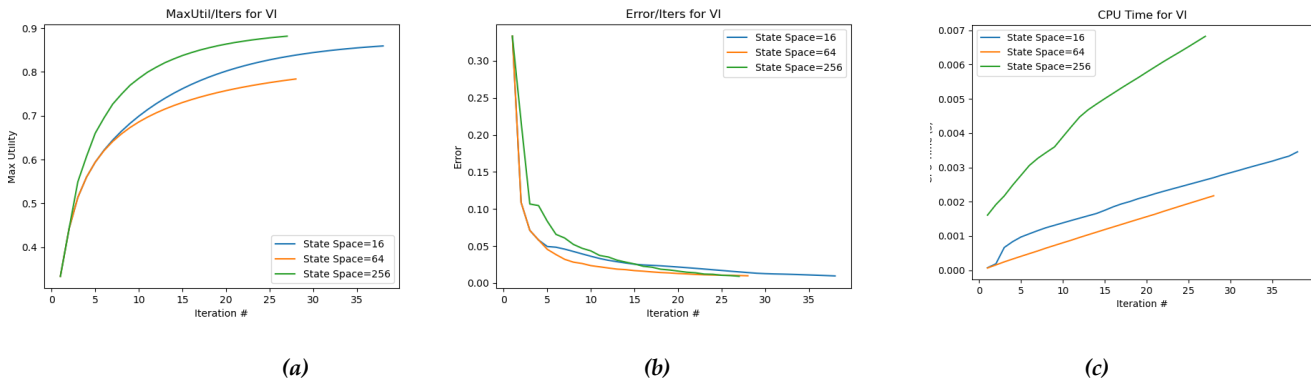


*(a)*     *(b)*     *(c)*

*Figure 6—(a)* FrozenLake, VI: Max Utility over Iterations for State Space
*(b)* FrozenLake, VI: Error over Iterations for State Space
*(c)* FrozenLake, VI: CPU Time over Iterations for State Space

# 4 EXPERIMENT 2: POLICY ITERATION

PI was used to evaluate the Forest Management problem. Figure 6a shows DiscountFactor=0.98 yielded the highest max utility at 53, emphasizing future rewards in determining the optimal policy (prioritizing 'Wait' occasionally as mentioned earlier) . It appears in 7a that the max utility values are unchanging over iterations for each discount value, but in fact there are slight fluctuations. While the max utility value is relatively constant, the number of iterations progresses until PI converges. **Error** is used to show that PI converges at 16 iterations for DiscountFactor=0.98 indicated by 7b, where in this case **error** is computed by taking the max delta between value functions associated with the current chosen policy and the next policy, i.e. the algorithm terminates when the policy is no longer changing.



(a)                                                  (b)

*Figure 7—(a)* Forest, PI: Max Utility over Iterations for Discount Factor
*(b)* Forest, PI: Error over Iterations for Discount Factor

For PI, there was no effect from toggling the number of states, as was the case for VI, so it is not shown here. The matrix representation of the optimal policy for Forest Management is shown in Figure 8 below.



*Figure 8—*Forest, PI: Visualization of Policy in Tuple Format

The policy indicates that a 'Wait' in State 0, followed by several 'Cut', and finally a series of 'Wait' yields the optimal value function. This is inline with the reward function that the problem defines, with rewards placed on 'Cut' and a 'Wait' or 'Cut' issued in the oldest state.

Next, PI was used on the Frozen Lake 4x4 problem, in which DiscountFactor=0.98 again yielded the highest max utility at 0.9 shown in 9a, indicating a strong emphasis on future rewards (reaching the goal ultimately). The algorithm converged at 5 iterations for DiscountFactor=0.98 as shown in 9b as the policy has no more variation. Other discount factor values take longer to converge, it seems, due to the policy continuing to change.
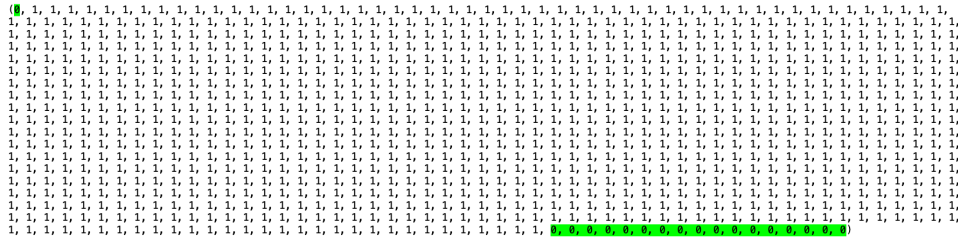


*(a)*                    *(b)*                    *(c)*

*Figure 9—(a)* FrozenLake, PI: Max Utility over Iterations for Discount Factor
*(b)* FrozenLake, PI: Error over Iterations for Discount Factor
*(c)* FrozenLake, PI: Visualization of Policy in Tuple Format

The visualization of the optimal policy for FrozenLake-4x4 grid is shown in 9c above. This graphic indicates that the PI agent has learned pathways to avoid holes on the way to the final goal, as all pathways point away from the existing holes on the grid.

## 5 PI AND VI ALGORITHM COMPARISON

Comparing VI and PI for the Forest problem in Figure 10, we see VI requires more iterations to converge, taking 49 iterations over PI's 17 iterations in 10a. While PI converges faster in terms of number of iterations, VI takes less CPU time to converge, with 0.08s vs 0.78s in 10b.



*(a)*                    *(b)*

*Figure 10—(a)* Convergence Comparison by Algorithm/MDP Problem
*(b)* CPU Time Comparison by Algorithm/MDP Problem

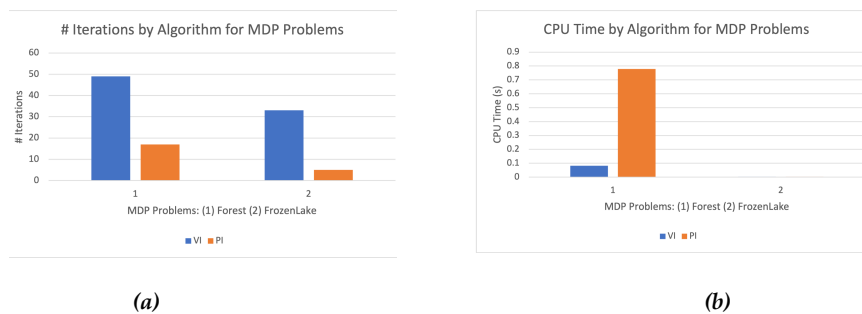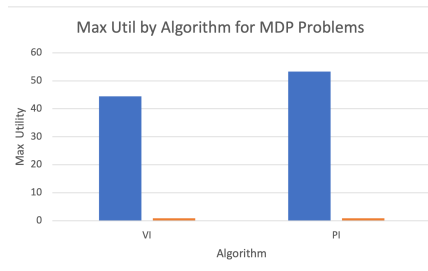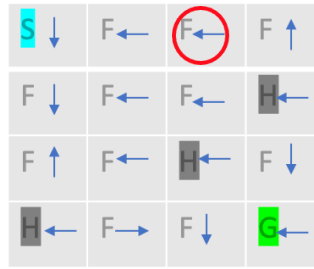It is logical that PI converges more quickly because a policy tends to converge in fewer iterations than a value function. This is because optimizing the *policy* first, instead of the value function, allows us to bypass a large set of potential value functions in the process of updating our policy. VI and PI are compared for FrozenLake-4x4 in Figure 10 above, with PI converging more quickly than VI as it did for the Forest problem. The same reasoning applies here.

To refresh from above, VI convergence was defined as having an error value below 0.01, where error is computed by taking the max delta between the previous and current value functions. PI convergence was defined to be when the policy is no longer changing.



*(a)*                          *(b)*                          *(c)*



*(d)*                          *(e)*

*Figure 11—(a)* Max Util Comparison by Algorithm/MDP Problem

*(b)* FrozenLake, VI: Visualization of Policy in Tuple Format

*(c)* FrozenLake, PI: Visualization of Policy in Tuple Format

*(d)* Forest, VI: Visualization of Policy in Tuple Format

*(e)* Forest, PI: Visualization of Policy in Tuple Format

For the Forest Problem, PI and VI did not converge to the same max utility (PI=53 and VI=44) as shown in Figure 11a, but their policies are very similar as shown in 11d and 11e above. The only difference between the two policies is that PI (11e) has one extra '0'/'Wait' action than VI (11d); despite the varying number of iterations and max utilities, both PI and VI have found *almost* identical policies. This is likely due to tuning the discount factor and providing the correct convergence criteria for both algorithms. The number of states did not affect anything for the Forest Management problem.

For the FrozenLake problem, it is noteworthy that VI and PI *do* converge to a very similar max utility (both around 0.9) as shown in 11a and very similar policies as well, shown in 11b and 10c. This is likely due to the problem space being small, and it being easier to find the optimal value function. The only difference between the two policies in 11b/11c is indicated by a red circle, which does not seem to affect the max utility significantly. The number of states *did* generate varied results for VI but not PI. This may be because even with a small problem space, finding the optimal value function can yield very different results. Finding the optimal policy is more likely to be identical across the small problems due to the limited set of state-action pairs.

## 6 EXPERIMENT 3: Q-LEARNING VS MODEL-BASED ALGORITHMS (VI/PI)

For the Forest problem, QLearner was trained for 400 episodes of 1000 iterations each, and hypertuned for Epsilon (probability of random action), DiscountFactor (importance on future rewards), Alpha (how large learning steps are), and EpsilonDecay (how much randomization decreases). Convergence was determined by finding where the policy was no longer changing beyond 400 episodes. Epsilon=0.75 suggests greater emphasis on random exploration of the space (12a) with a gradual EpsilonDecay=0.94 (12d), DiscountFactor=0.98 emphasizes future rewards (12b), and Alpha=0.8 signifies higher learning rate to traverse the large space (12c). Tuning Epsilon/Discount provide the most interesting information, with high values indicating this is *not* a greedy learner (it prioritizes exploring/ future rewards over short-term Q-value).



*Figure 12—(a)* Forest, QLearning: Max Utility over Iterations for Epsilon
*(b)* Forest, QLearning: Max Utility over Iterations for Discount Factor
*(c)* Forest, QLearning: Max Utility over Iterations for Learning Rate
*(d)* Forest, QLearning: Max Utility over Iterations for Epsilon Decay

Figure 13 shows QLearning has a smaller max utility and longer computation time compared to VI/PI. This could be because VI/PI are model-based, planning algorithms that know the model and reward function ahead of time. QLearning being a model-free algorithm only computes

*expected* rewards based on state-action pairs, and therefore requires more iterations to explore the space and may not ultimately find the optimal solution compared to VI/PI in this case.



*(a)*                                                                                   *(b)*

***Figure 13—(a)*** Max Util Comparison by Algorithm for Forest Management
***(b)*** CPU Time Comparison by Algorithm for Forest Management

Figure 14 below shows the optimal policy for this QLearner to better understand why we see PI/VI outperform QLearner. Compared to the tuples above for VI/PI, this policy's 0's and 1's are clearly more scattered. Given that the reward function assigns 1 point for 'Cut'/'1' in an intermediate state, and 4 points for 'Wait'/'0' in the final state, it makes sense that VI/PI performs better over QLearning. It is understandable that QLearner would not come to the same conclusion as VI/PI, given that it did not know that reward-scheme in advance.

```
(0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0,
0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1,
0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0,
1, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1,
0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1,
0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0,
1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0,
0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1,
1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0,
0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1,
1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1,
0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0,
1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1,
1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1,
0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0,
1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0,
0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0)
```

***Figure 14—*** Forest, QLearning: Visualization of Policy in Tuple Format

QLearner was used on the FrozenLake-4x4 problem as well, with 400 episodes of 1000 iterations each used for training. Epsilon=0.6 and DiscountFactor=0.98 yielded the highest max utility, as shown in Figure 15a and 15b below. This suggests moderate exploration is needed (reasonable as the problem space is small) and an emphasis placed on future rewards. Other parameter tuning did not result in any significant gains. The convergence criteria was determined the same way as above for the Forest problem.
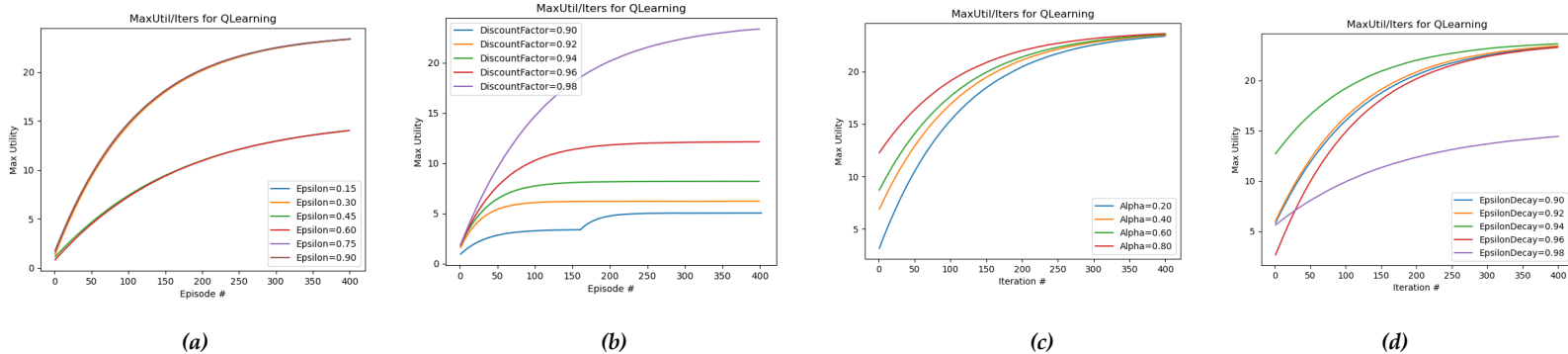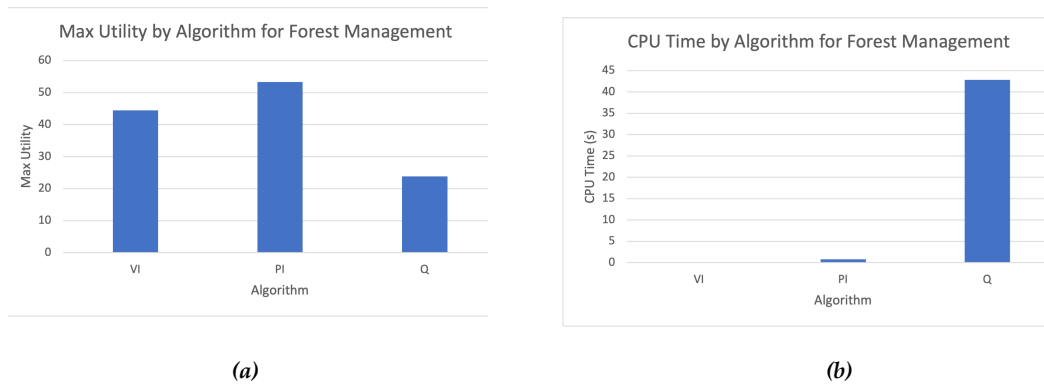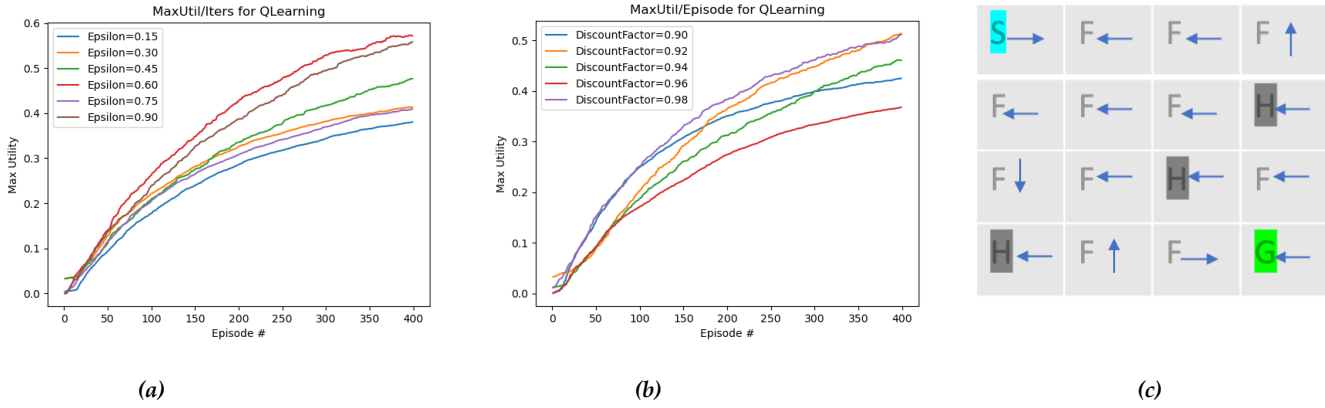
*Figure 15*—*(a)* FrozenLake, QLearning: Max Utility over Iterations for Epsilon

*(b)* FrozenLake, QLearning: Max Utility over Iterations for Discount Factor

*(c)* FrozenLake, QLearning: Visualization of Policy in Tuple Format

Figure 15c shows the optimal policy visualization for FrozenLake-4x4 grid. While there are overlapping components with PI/VI, a few states point *towards* the holes, which would cause the agent to get 0 reward. This explains the performance comparison between the algorithms.
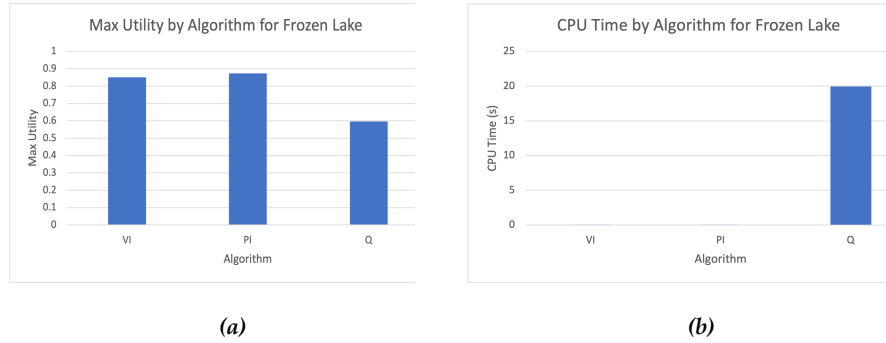


*Figure 16*—*(a)* Max Util Comparison by Algorithm for Frozen Lake

*(b)* CPU Time Comparison by Algorithm for Frozen Lake

The comparison of QLearning vs VI/PI in Figure 16 shows that QLearning slightly underperforms VI/PI in terms of max utility in 16a, and does take more computation time than VI/PI in 16b, whose compute times are negligible. Overall, QLearning seems to perform better on this smaller problem than the large Forest Management problem, as the state space is more manageable and the agent can compute expected rewards for the states that start to mimic the true reward function. The better performance of VI/PI can be explained by the fact that VI/PI are model-based, and therefore have the model and reward/transition functions beforehand, which requires less exploration than for QLearning.