# Reinforcement Learning - Week 2

Antonio Chiappetta

October 2019

## 1 Sailor gridworld

### 1.1 Question 1

The agent is the sailor. He is the one that has to make decisions to perform actions, taking into consideration the environment around him and the goal he wants to reach, i.e., the harbour.

The environment is everything that is outside of the agent's absolute control: the sea, the rocks, the wind, the harbour, and the way these elements influence the outcome of the actions taken by the sailor.

## 2 Value iteration

### 2.1 Task 1

The following image shows the final render after 100 iterations of the value iteration algorithm. In this image only state values are computed, not the policy.
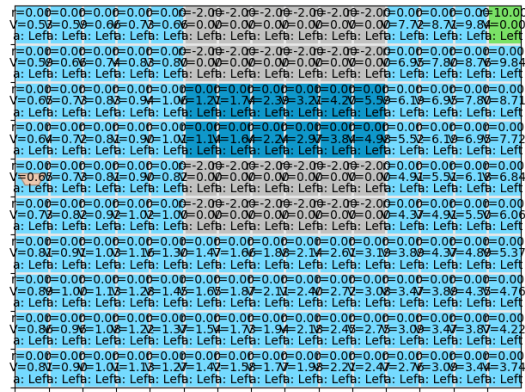


Figure 1: State values after 100 value iterations

## 2.2 Question 2

The state value of the harbour and rock states is 0. Both are termination states, the episode finishes in both cases: in the former case because the objective has been reached, and in the latter because the sailor dies of hitting the rock. The termination implies there are no other rewards following those states, so the discounted return provided in that location is obviously null.

## 2.3 Task 2

The following images shows the final render in 3 executions of the program where 100 iterations of the value iteration algorithm were used both to update state values and to determine an optimal policy. The sailor always managed to find his way to the goal.
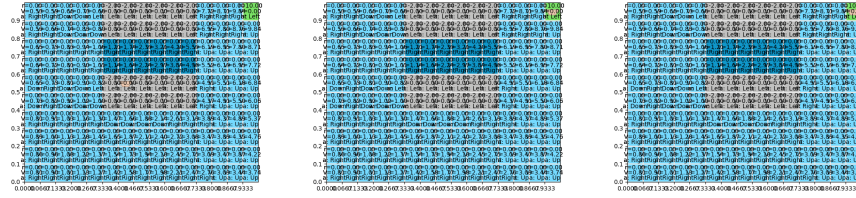
Figure 2: State values after 100 value iterations

## 2.4 Question 3

During the previous executions, the sailor was always trying to use the short path between the rocks to reach the goal. Despite the possible penalty resulting in hitting the rocks, the great difference in distance to the goal compared to the path going around the goal made this solution still optimal.

But if we change the penalty for hitting the rocks to -10, making the sailor value his life more, he will absolutely try to avoid being closer to the rocks and will choose a different path, going towards the bottom of the grid and then right towards the column of the harbour. The behaviour we saw showed that the sailor was always keeping himself at least 2 steps away from the rocks.
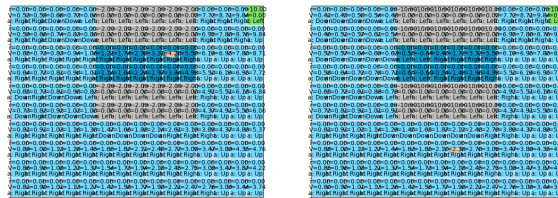
Figure 3: Different behaviour depending on rocks penalty

2

## 2.5 Question 4

I tried running the algorithm for 50, 25, 20 and 10 iterations. The final grid with value and policy functions is shown in the next figures for all these cases.
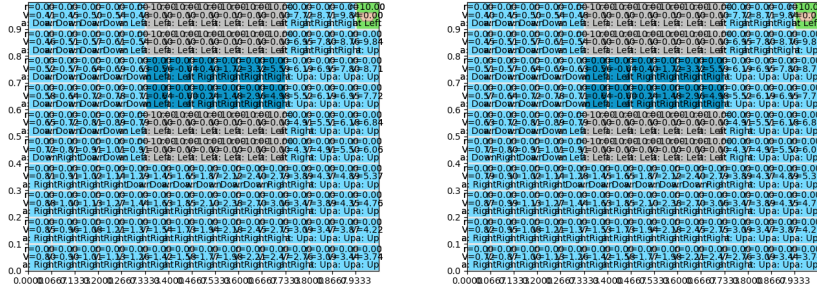


Figure 4: Value and policy functions convergence with 50 and 25 iterations
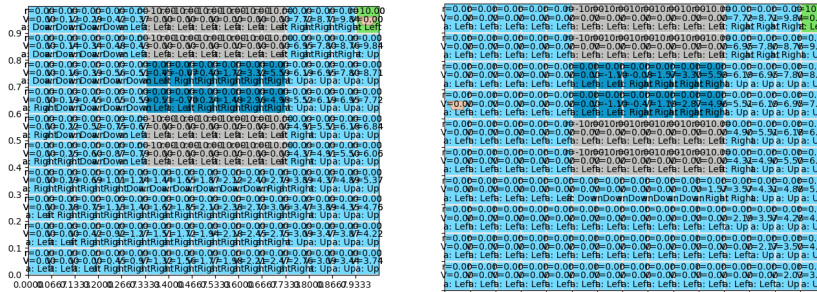


Figure 5: Value and policy functions convergence with 20 and 10 iterations

It's immediate to notice that the policy function converges before the value function: once the relative order among the state values is fixed, the optimal policy has been found; from that moment on, state values will keep changing and adjusting to the optimal values, but the policy will stay the same.

With 50 iterations, both functions have already converged as we can see from the state and action values that are the same as in the last figure of the previous exercise, since we ran the computation again with a penalty of -10 for hitting the rocks .

With 25 iterations the value function has slightly different values, meaning it did not converge yet to the values of 50 iterations (at least up to 2 decimal digits). But the relative order of states it's the same and we can see that the policy function has converged.

With 20 iterations the policy function has not converged yet but the sailor is still able to find his way to the goal.

With 10 iterations the sailor is not even able to find a path to reach the harbour.

## 2.6   Task 3

Changing the termination condition with the use of a small threshold, a new run of the algorithm (with rock penalty back to -2) terminated the value iteration in 31 steps. The following figure shows the results of an episode after this execution.



Figure 6: Value iteration with threshold as termination condition

## 2.7   Task 4

Running the program for $N = 1000$ episodes, the discounted return of the initial state has a mean $\mu(R_0) = 0.67$ and a standard deviation $\sigma(R_0) = 1.36$.

## 2.8   Question 5

Since we are executing the optimal policy, the result regarding the discounted return is coherent with the optimal value function calculated, that says that the value of the initial state is $V^*(s_0) = 0.66$. The value function tells the value of a state in terms of expected return when starting in that state and following the chosen policy thereafter; that's why, over a great number of episodes, the average discounted return is expected to tend to this value.

## 2.9   Question 6

The value iteration approach used in the previous tasks assumes complete knowledge of the environment, i.e., having the complete probability distributions of all

possible transitions between states. Given a state and an action, with this information it's possible to generate predictions of the resultant next state and next reward, thus simulating the environment to calculate value and policy functions. [1]

However, in an unknown environment, this information is not available and real experience generated by the environment is needed to learn these probability distributions.

# References

[1] R. S. Sutton and G. Barto. *Reinforcement Learning*. The MIT Press, 2015.