

# RL-2

Abhilash Jain

October 13, 2019

## 1 Sailor Grid World

### 1.1 What is the agent and the environment in this setup?

The Agent is the sailor who wants to reach the harbour.

The environment is the sea, the rocks, the harbour and the narrow passage between the rocks. Also the calmness of the sea and the heavy windy conditions are part of the environment

## 2 Value iteration

### 2.1 Task 1

```
1 def action_v(state, V):
2     A = []
3     for transition in env.transitions[state[0], state[1]]:
4         action_value = 0
5         for next_state, reward, done, prob in transition:
6             action_value += prob * (reward + (0.9 *
7                                     V[next_state] if not done else
8                                     0))
9         A.append(action_value)
10    return A
11
12 # TODO: Compute the value function and policy (Tasks 1, 2 and 3)
13 value_est, policy = np.zeros((env.w, env.h)), np.zeros((env.w,
14 env.h))
15 theta = value_est.copy()
16 eps = 0.0001
17 width,height=env.w,env.h
18 for itera in range(100):
19     env.clear_text()
20     delta=0
21     for i in range(width):
22         for j in range(height):
23             A = action_v((i, j), value_est)
24             value_est[i,j] = np.max(A)
25             policy[i,j] = np.argmax(A)
26     env.draw_values(value_est)
27     env.draw_actions(policy)
```

Listing 1: Task 1

State Values:

	r=0.00 V=0.53 a: Right	r=0.00 V=0.59 a: Right	r=0.00 V=0.66 a: Down	r=0.00 V=0.73 a: Down	r=0.00 V=0.66 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=0.00 V=7.72 a: Right	r=0.00 V=8.71 a: Right	r=0.00 V=9.84 a: Right	r=10.00 V=0.00 a: Left	
0.9	r=0.00 V=0.59 a: Right	r=0.00 V=0.66 a: Right	r=0.00 V=0.74 a: Down	r=0.00 V=0.83 a: Down	r=0.00 V=0.80 a: Down	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=0.00 V=6.95 a: Right	r=0.00 V=7.80 a: Right	r=0.00 V=8.76 a: Right	r=0.00 V=9.84 a: Up	
0.8	r=0.00 V=0.65 a: Right	r=0.00 V=0.73 a: Right	r=0.00 V=0.83 a: Right	r=0.00 V=0.94 a: Right	r=0.00 V=1.06 a: Right	r=0.00 V=1.21 a: Right	r=0.00 V=1.74 a: Right	r=0.00 V=2.39 a: Right	r=0.00 V=3.21 a: Right	r=0.00 V=4.20 a: Right	r=0.00 V=5.59 a: Right	r=0.00 V=6.19 a: Up	r=0.00 V=6.95 a: Up	r=0.00 V=7.80 a: Up	
0.7	r=0.00 V=0.64 a: Right	r=0.00 V=0.72 a: Right	r=0.00 V=0.80 a: Right	r=0.00 V=0.90 a: Right	r=0.00 V=1.01 a: Right	r=0.00 V=1.14 a: Right	r=0.00 V=1.64 a: Right	r=0.00 V=2.24 a: Right	r=0.00 V=2.97 a: Right	r=0.00 V=3.84 a: Right	r=0.00 V=4.98 a: Right	r=0.00 V=5.52 a: Up	r=0.00 V=6.19 a: Up	r=0.00 V=6.95 a: Up	
0.6	r=0.00 V=0.65 a: Right	r=0.00 V=0.73 a: Right	r=0.00 V=0.81 a: Down	r=0.00 V=0.90 a: Down	r=0.00 V=0.82 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=0.00 V=4.91 a: Right	r=0.00 V=5.51 a: Up	r=0.00 V=6.18 a: Up	r=0.00 V=6.84 a: Up	
0.5	r=0.00 V=0.72 a: Down	r=0.00 V=0.81 a: Right	r=0.00 V=0.91 a: Down	r=0.00 V=1.02 a: Down	r=0.00 V=1.00 a: Down	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Left	r=0.00 V=4.37 a: Right	r=0.00 V=4.91 a: Up	r=0.00 V=5.50 a: Up	r=0.00 V=6.06 a: Up	
0.4	r=0.00 V=0.80 a: Right	r=0.00 V=0.91 a: Right	r=0.00 V=1.03 a: Right	r=0.00 V=1.16 a: Right	r=0.00 V=1.30 a: Right	r=0.00 V=1.47 a: Down	r=0.00 V=1.66 a: Down	r=0.00 V=1.88 a: Down	r=0.00 V=2.14 a: Down	r=0.00 V=2.61 a: Right	r=0.00 V=3.19 a: Right	r=0.00 V=3.89 a: Right	r=0.00 V=4.37 a: Up	r=0.00 V=4.89 a: Up	
0.3	r=0.00 V=0.88 a: Right	r=0.00 V=1.00 a: Right	r=0.00 V=1.13 a: Right	r=0.00 V=1.28 a: Right	r=0.00 V=1.45 a: Right	r=0.00 V=1.65 a: Right	r=0.00 V=1.87 a: Right	r=0.00 V=2.11 a: Right	r=0.00 V=2.40 a: Right	r=0.00 V=2.72 a: Right	r=0.00 V=3.08 a: Right	r=0.00 V=3.47 a: Up	r=0.00 V=3.89 a: Up	r=0.00 V=4.35 a: Up	
0.2	r=0.00 V=0.83 a: Right	r=0.00 V=0.95 a: Right	r=0.00 V=1.08 a: Right	r=0.00 V=1.22 a: Right	r=0.00 V=1.37 a: Right	r=0.00 V=1.54 a: Right	r=0.00 V=1.73 a: Right	r=0.00 V=1.94 a: Right	r=0.00 V=2.18 a: Right	r=0.00 V=2.45 a: Right	r=0.00 V=2.75 a: Right	r=0.00 V=3.09 a: Up	r=0.00 V=3.47 a: Up	r=0.00 V=3.87 a: Up	
0.1	r=0.00 V=0.72 a: Right	r=0.00 V=0.88 a: Right	r=0.00 V=1.00 a: Right	r=0.00 V=1.13 a: Right	r=0.00 V=1.27 a: Right	r=0.00 V=1.42 a: Right	r=0.00 V=1.58 a: Right	r=0.00 V=1.77 a: Right	r=0.00 V=1.98 a: Right	r=0.00 V=2.21 a: Right	r=0.00 V=2.47 a: Right	r=0.00 V=2.76 a: Right	r=0.00 V=3.09 a: Up	r=0.00 V=3.44 a: Up	
0.0	r=0.00 V=0.72 a: Right	r=0.00 V=0.88 a: Right	r=0.00 V=1.00 a: Right	r=0.00 V=1.13 a: Right	r=0.00 V=1.27 a: Right	r=0.00 V=1.42 a: Right	r=0.00 V=1.58 a: Right	r=0.00 V=1.77 a: Right	r=0.00 V=1.98 a: Right	r=0.00 V=2.21 a: Right	r=0.00 V=2.47 a: Right	r=0.00 V=2.76 a: Right	r=0.00 V=3.09 a: Up	r=0.00 V=3.44 a: Up	
	0.0000	0.0667	0.1333	0.2000	0.2667	0.3333	0.4000	0.4667	0.5333	0.6000	0.6667	0.7333	0.8000	0.8667	0.9333

## 2.2 Question 2

What is the state value of the harbour and rock states? Why?

The State Value of the harbour and states as seen is zero, as they are the termination conditions or terminal states. The value of a state is the expected sum of all future rewards when starting in that state and following a specific policy. For the terminal state, this is zero - there are no more rewards to be had.

## 2.3 Task 2

```
1     done=False
2     while not done:
3         # Select a random action
4         # TODO: Use the policy to take the optimal action (Task
5         2)
6         #action = int(np.random.random()*4)
7         action=policy[state]
8         # Step the environment
9         state, reward, done, _ = env.step(action)
10        # Render and sleep
11        env.render()
        sleep(0.1)
```

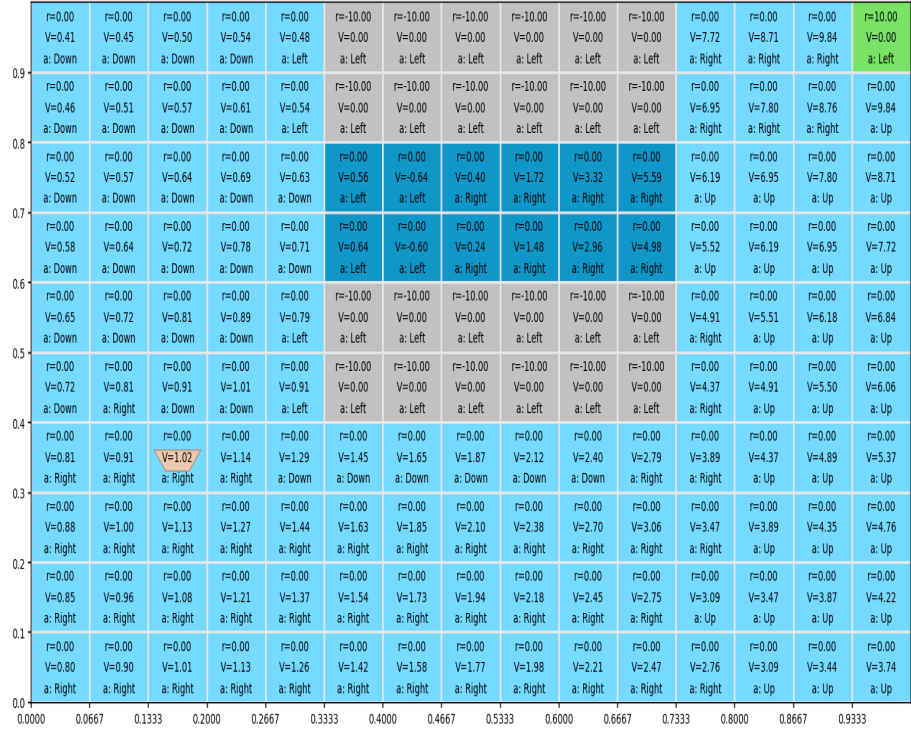
Listing 2: Task 2

The sailor is able to reach the harbour most of the times using the optimal path, whenever it failed, it was crashing into the rocks

## 2.4 Question 3

Which path did the sailor choose? If you change the reward for hitting the rocks to -10 (that is, make the sailor value his life more), does he still choose the same path?

The sailor takes the longer path instead of the narrower path as it is a more safer option. As you can see the Value function is very low for the states in the beginning of the narrow package and the other path has higher values.



## 2.5 Question 4

What happens if you run the algorithm for a smaller amount of iterations? Do the value function and policy still converge? Which of them - the policy or value function - needs less iterations to converge, if any? Justify your answer.

The algorithm at-least takes 30 iterations to converge, anything below that it doesn't converge for example it didn't converge at 15 iterations. Both the policy and the value function take the same time to converge, at-least that is what i observed. It is due to the fact that, as the Values keep changing, so does the policy, and once the Values have converged so has the policies.

This was checked by doing a simple comparison(`numpy.array_equal`), and both the policies converged at the 30th epoch and the value function converged as well at the 30th epoch (given the small epsilon)

But something also observed was, the optimal policy(shortest path from the narrow passage) was learnt faster than the value function i.e the values in the optimal path were still changing, but the policies werent.

## 2.6 Task 3

Set the reward for crashing into the rocks back to -2. Change the termination condition of your algorithm to make it run until convergence. You can assume the values to have converged if the maximum change in value is lower than a certain threshold :

```
1 # TODO: Compute the value function and policy (Tasks 1, 2 and 3)
2 value_est, policy = np.zeros((env.w, env.h)), np.zeros((env.w,
3 env.h))
4 theta = value_est.copy()
5 eps = 0.0001
6 width,height=env.w,env.h
7 for itera in range(100):
8     env.clear_text()
9     delta=0
10    for i in range(width):
11        for j in range(height):
12            A = action_v((i, j), value_est)
13            value_est[i,j] = np.max(A)
14            policy[i,j] = np.argmax(A)
15            #From here epsilon code impleted
16            delta=max(delta,np.abs(value_est[i,j]-theta[i,j]))
17    if (delta < eps):
18        print('Stopped',itera)
19        break
20    else :
21        theta = np.copy(value_est)
22        #To check values every iteration
23        #env.draw_values(value_est)
24        #env.draw_actions(policy)
25        #env.render()
26
27    # Show the values and the policy
28    env.draw_values(value_est)
29    env.draw_actions(policy)
30    env.render()
31    sleep(1)
```

Listing 3: Task 3

## 2.7 Task 4

Modify the program to compute the discounted return of the initial state. Create a loop to run the program for N 1000 episodes and see what values you get. Report the average and standard deviation.

```
1 # TODO: Run multiple episodes and compute the discounted returns (
2 Task 4)
3 done = False
4 discount_r=[]
5 for i in range(1000):
6     state = env.reset()
7     temp_discount_r=0
8     done=False
9     j=0
```

```

9         while not done:
10             # Select a random action
11             # TODO: Use the policy to take the optimal action (Task
12             2)
13             action = int(np.random.random()*4)
14             action=policy[state]
15             # Step the environment
16             state, reward, done, _ = env.step(action)
17             temp_discount_r+=(0.9**j) * reward
18             #Render and sleep
19             #env.render()
20             #sleep(0.1)\
21             j+=1
22             discount_r.append(temp_discount_r)
23         print("Mean",mean(discount_r))
24         print("std deviation",np.std(np.array(discount_r)))

```

Listing 4: Task 3

Mean 0.6578227696413183

std deviation 1.3761068012248063

Which is very close to the Value function of the state the sailor started in

## 2.8 Question 5

What is the relationship between the discounted return and the value function? Explain briefly.

As we can clearly observe that the mean of discounted rewards over 1000 episodes is very close to the value function of the state the sailor started in , we can say that the value function is the expectation of the discounted reward. As we do more and more iterations both the values will eventually converge

Cross checking this with setting the state of the sailor to start in 13,9 the value function is 9.44 and the average obtained is 9.4335466 which is again very close

## 2.9 Question 6

Imagine a reinforcement learning problem involving a robot exploring an unknown environment. Could the value iteration approach used here be applied directly to that problem? Why/why not? Which of the assumptions are unrealistic, if any?

No the value iteration cannot be applied directly in the robot exploring an unknown environment because

- Value iteration assumes that we know all the states, including the end states, and is a finite state space but in an unknown environment that is not possible
- In an unknown environment we also don't know the end state from the start, it would not be possible to iterate over a space for value iteration

when we don't know the "end state", as Value iteration starts from the end state, and works it way backwards

- The transition probabilities between states are not available in an unknown environment, making it not possible not 'iterate' over and we wont be able to define a policy