

Exercise 2

Gengcong Yan - 1009903
ELEC-E8125 - Reinforcement Learning

October 10, 2021

1 Question 1

- Agent: The sailor
- Environment: Water area, Harbour, Narrow passage, Rocks, Possible wind.

2 Task 1

```
1 def values_iterations(env, gamma=1.0, itNum=100):
2
3     value_est, policy = np.zeros((env.w, env.h)), np.zeros((env
4         .w, env.h))
5
6     for it in range(itNum):
7         # iterate every state in 2D
8         env.clear_text()
9         for w in range(env.w):
10             for h in range(env.h):
11                 Q_value=[]
12                 for tran in env.transitions[w,h]:
13                     # print("tran:")
14                     # print(tran)
15                     A_value=0
16                     for next_state, reward, done, prob in tran:
17                         # consider all situation under
18                         # different prob
19                         A_value+=prob*(reward+(gamma*value_est[
20                             next_state] if not done else 0))
21                     Q_value.append(A_value)
22                     # print(it, Q_value)
23                 value_est[w,h]=np.max(Q_value)
24                 policy[w,h]=np.argmax(Q_value)
25
26     return value_est, policy
```

The render environment of Task 1 is shown in Fig.1.

3 Question 2

The reward value of harbour and rock states are 10 and -2, respectively. But the state values of them are both 0. Because in our default setup of environment, the episode terminates as soon as the sailor arrives harbour and rocks. There are no next state for them to transfer, their values can't get updated then always be zero.

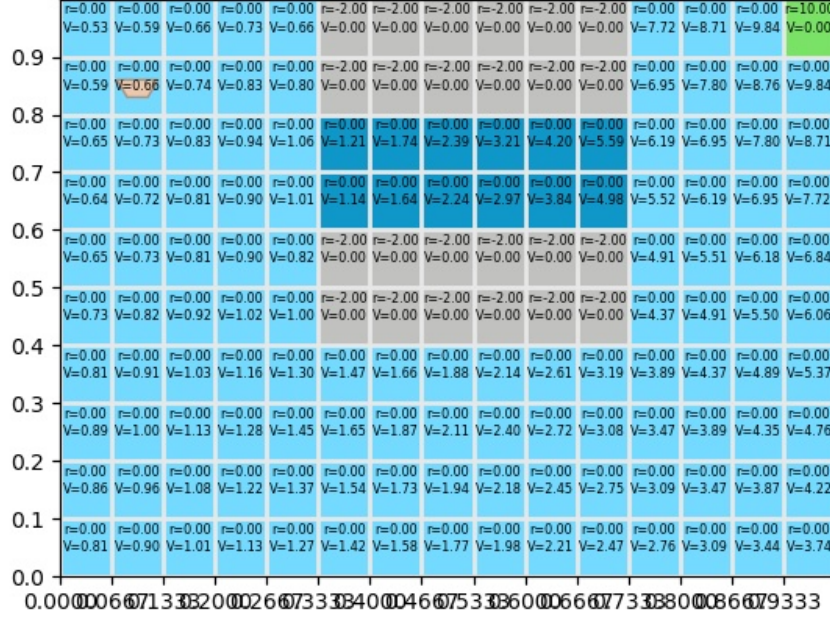


Figure 1: Values in environment after 100 iterations

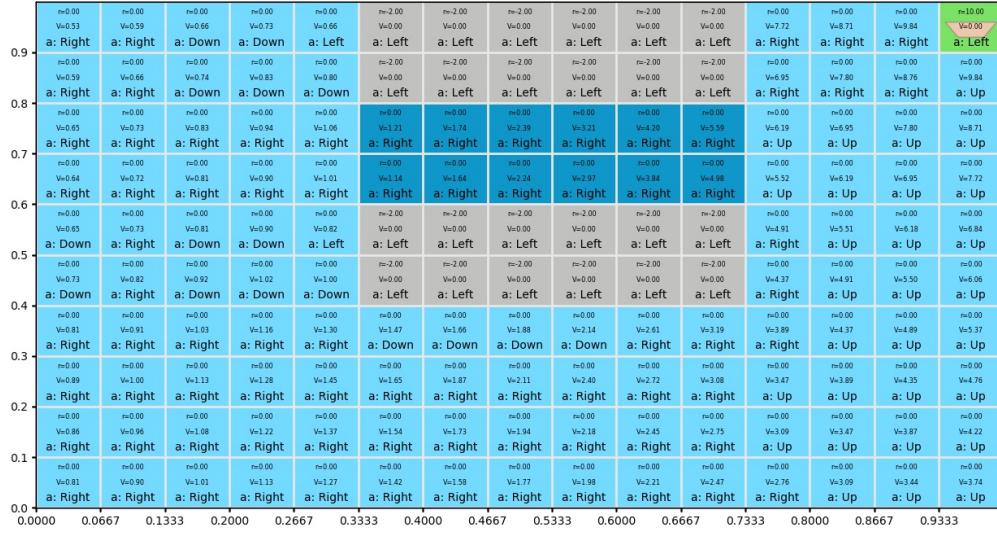


Figure 2: Actions in environment after 100 iterations

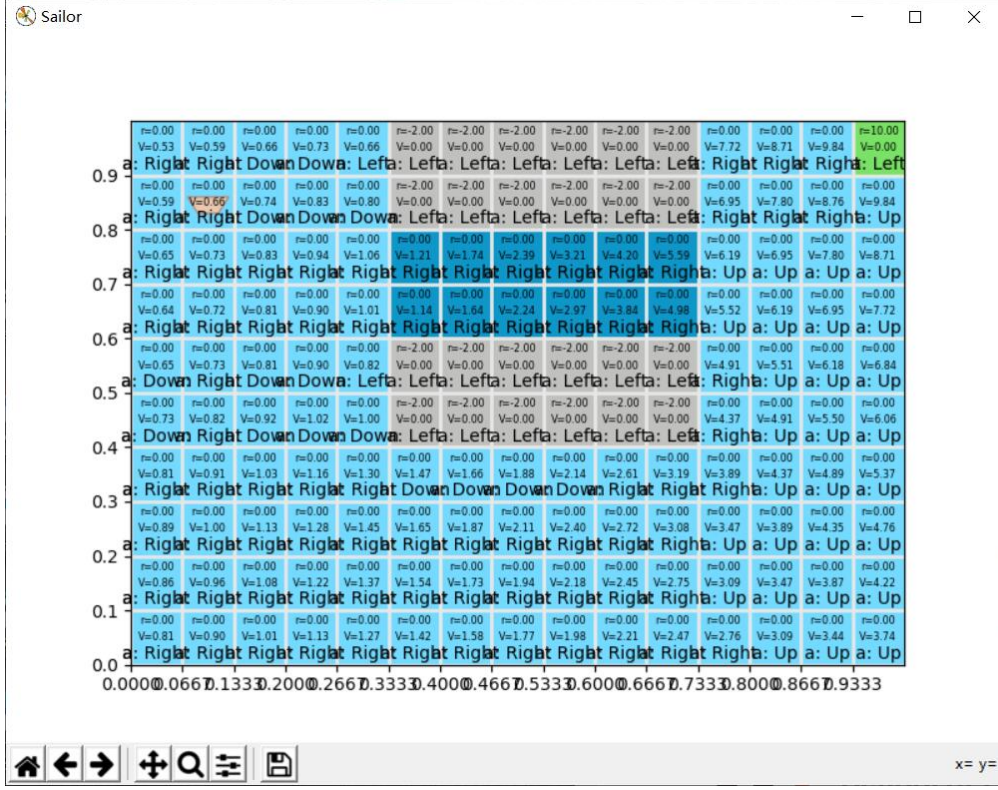
4 Task 2

The actions in environment is shown in Fig.2.

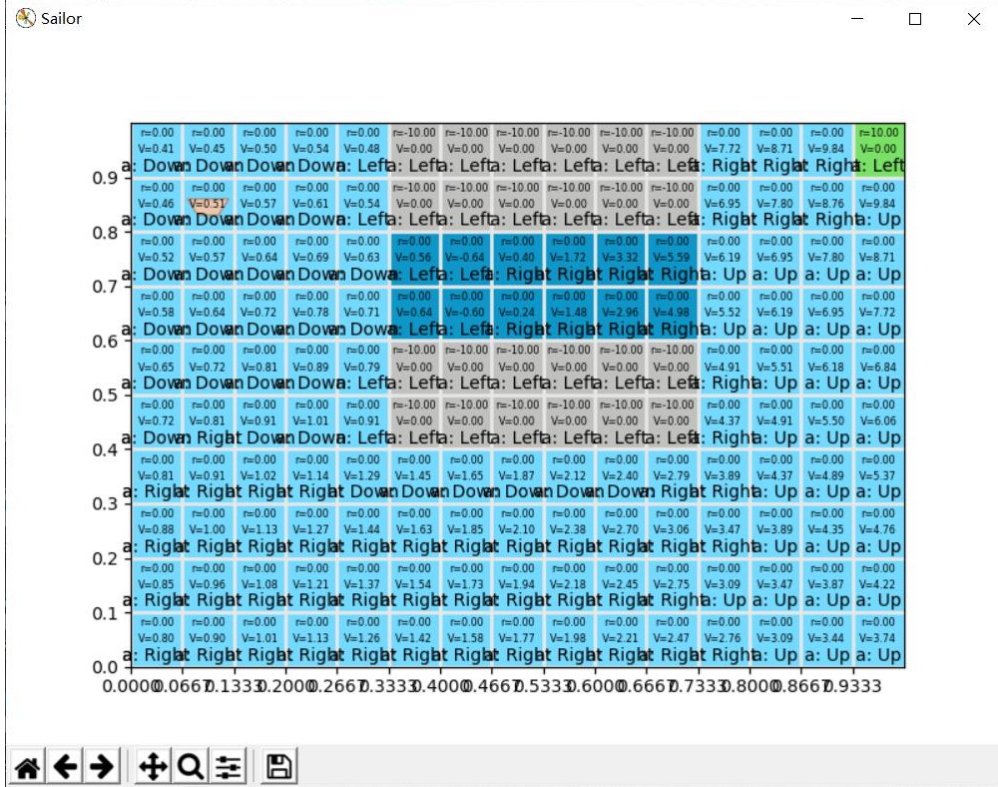
5 Question 3

- Rock reward=-2. The sailor tend to risk choosing the dangerous path between the rocks.
- Rock reward=-10. The sailor will choose the safe path below the rocks due to the smaller negative reward. The sailor value their own life more.

The comparison is in Fig.3.



(a) Rock reward:-2



(b) Rock reward:-10

Figure 3: Choices in different rock rewards

6 Question 4

```
1 def values_iterations(env,gamma=1.0,itNum=100):
2     value_est, policy = np.zeros((env.w, env.h)), np.zeros((env
3         .w, env.h))
4
5     old_values=value_est.copy()
6     old_policy=policy.copy()
7
8     for it in range(itNum):
9         # iterate every state in 2D
10        env.clear_text()
11        for w in range(env.w):
12            for h in range(env.h):
13                Q_value=[]
14                for tran in env.transitions[w,h]:
15                    A_value=0
16                    for next_state,reward,done,prob in tran:
17                        # consider all situation under
18                        # different prob
19                        A_value+=prob*(reward+(gamma*value_est[
20                            next_state] if not done else 0))
21                    Q_value.append(A_value)
22                # print(it,Q_value)
23                value_est[w,h]=np.max(Q_value)
24                policy[w,h]=np.argmax(Q_value)
25
26        # value difference checks converge
27        if ((value_est-old_values)<epsilon).all():
28            print("converge at "+ str(it))
29            break
30        else:
31            # print("Not converge at "+ str(it))
32            old_values=np.copy(value_est)          # should
33            use np.copy, directly equal causes fault
34
35        # policy difference checks converge
36        if ((policy-old_policy)<epsilon).all():
37            print("Policy converge at "+ str(it))
38            break
39        else:
40            print("Policy didn't converge at "+ str(it))
41            old_policy=np.copy(policy)
42    return value_est,policy
```

I implemented the convergence check in original function at Task 1, using the formula in introduction file. There are two types of data to detect convergence, the value function and

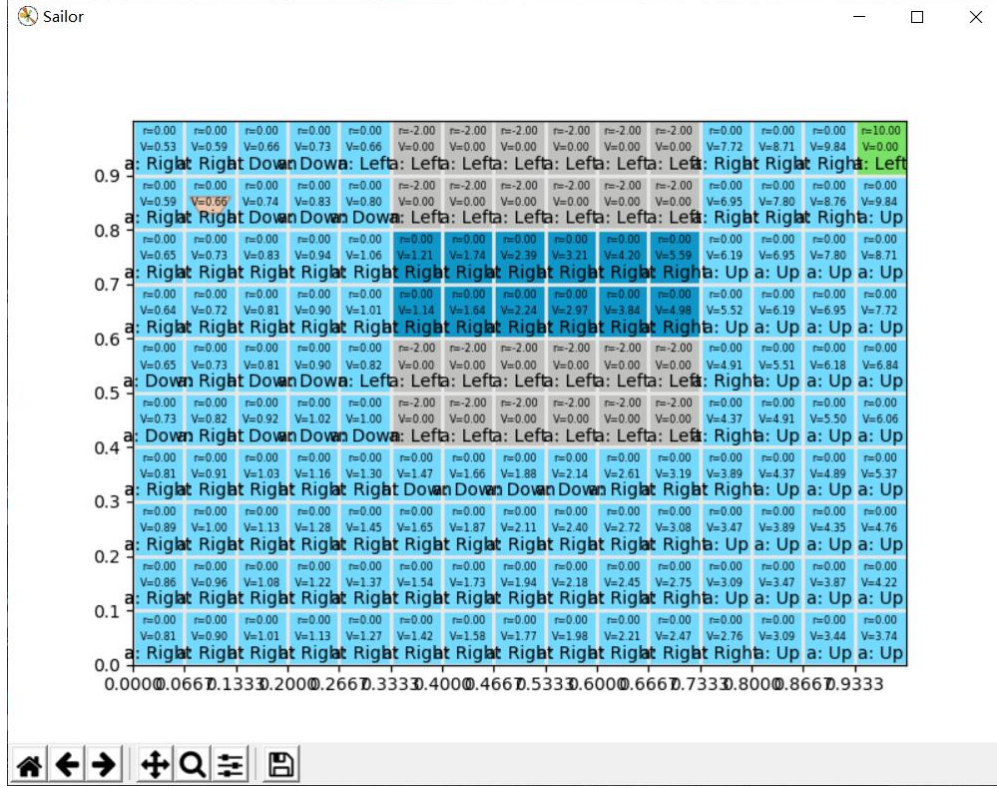


Figure 4: Convergence at rock reward=-2

the policy. After test in both of them, it shows they converged at different iterations. from my experiment, the value function converged at iteration 29, whereas the policy at 23. The policy converged sooner than the value function does. The reason is that the value function needs the whole state space to take into account. A small change in state values does not necessarily lead to a change in policy.

7 Task 3

The algorithm converged at iteration 39, when the rock reward equals -2. The grid is shown in Fig.4.

8 Task 4

After 1000 episodes:

- Average: 0.683
- Standard deviation: 1.361

```

1     results=[]
2     for ep in range(1000):
3         discount_return=0
4         i=0
5         done = False
6         while not done:
7             # Select a random action
8             # TODO: Use the policy to take the optimal action (
              Task 2)
9             action = policy[state]
10
11            # Step the environment
12            state, reward, done, _ = env.step(action)
13            discount_return+=reward*(gamma**i)
14            i+=1
15            # Render and sleep
16            # env.render()
17            # sleep(0.2)
18
19            print("Episode"+ str(ep)+" : "+ str(round(
              discount_return,5)))
20            results.append(discount_return)
21
22            state = env.reset()
23
24            # print(results)
25            print("-----Average& STD") # 0.683 & 1.361
26            print(np.average(results))
27            print(np.std(results))

```

9 Question 5

From the equation below we know, the value function of a state s under a policy π , denoted $v(s)$, is the expected discounted return when starting in s and following thereafter. [Chapter 3 Sutton & Barto]

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s] \quad (1)$$

10 Question 6

No. It cannot directly apply into the problem involving a robot exploring an **unknown** environment. Mainly because in our exercise, we have finite state space of environment with certain limited number of actions. What's more, in a more real world setting computing the values in different states like we did in this exercises will be computationally expensive. You can

see there are 3 or 4 loops in code, requiring a lot of time and computing power when states become numerous. Lastly, the various states and actions in exercise have been abstractly expressed so clearly while it is difficult to abstract the complexity of the situation in a real application problem.

11 Feedback

1. 5h