

Exercise 1

Breno Aberle - 876438
ELEC-E8125 - Reinforcement Learning

September 28, 2020

1 Question 1

In that Reinforcement Learning exercise the agent is the sailor. The sailor wants to reach the harbour because there is a positive reward. The agent learns with every episode. Given the current state it has to take an action. The agent can take four actions: left, down, right, up. The wind probability, effects of wind, and exact location of the harbour is perfectly known to the agent.

The environment is the whole cell grid with harbour, rocks, water with less wind and water with more wind. The harbour has a reward of +10, rocks -2 and episode termination and the rest a reward of 0. The sea has a low wind probability, while the probability of wind in the narrow passage between the rocks higher is.

2 Task 1

The following code represents the value iteration:

```

1  gamma = 0.9 # discount factor value
2  number_iterations = 100
3  actions = [0,1,2,3]
4  value_est_history = []
5  policy_history = []
6
7  for i in range(number_iterations):
8      env.clear_text() # remove previously drawn values
9      value_est_copy = value_est.copy()
10     policy_copy = policy.copy()
11     value_est_history.append(value_est_copy)
12     policy_history.append(policy_copy)
13
14     for x in range(env.w): # loop through all states/cell in
        the environment
15         for y in range(env.h):
16             state_values_of_actions = [] # keep track of
                calculated state values of each action
17             for action in actions:
18                 transitions = env.transitions[x,y,action]
19                 state_value_of_action = 0
20                 for transition in transitions:
21                     state_next = transition[0]
22                     reward = transition[1]
23                     done = transition[2]
24                     probability = transition[3]
25                     if (state_next == None):
26                         state_value_of_action += 0
27                     else:
28                         state_value_of_action += probability *
                            (reward + gamma * value_est_copy[
                                state_next])
29                 state_values_of_actions.append(
                    state_value_of_action)
30             value_est[x][y] = np.max(state_values_of_actions)
31             policy[x][y] = np.argmax(state_values_of_actions)
32
33     max_diff_val = np.max(abs(abs(value_est) - abs(
        value_est_copy)))
34     if (max_diff_val < epsilon):
35         print("Converged! Value state converged in iteration: "
            , i+1)
36         break
37     max_diff_policy = np.max(abs(policy_copy - policy))
38     if (max_diff_policy < epsilon):
39         print("Converged! Policy converged in iteration: ", i
            +1)

```

The figure 1 shows an example render after convergence.

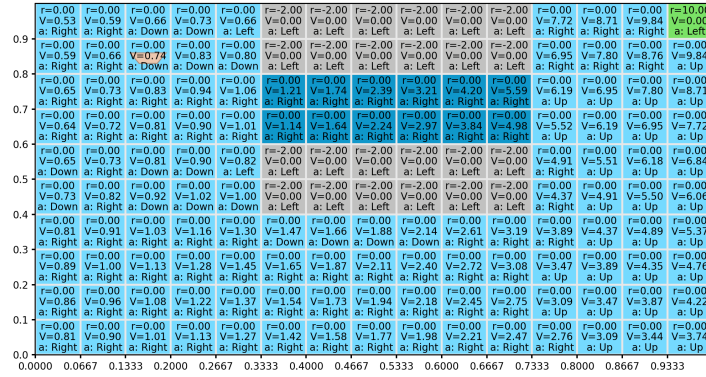


Figure 1: Environment with reward, state value, and policy for each cell.

3 Question 2

The state value of the harbour as well as the rocks is 0. That's because after one of these are reached the episode gets terminated. So there won't follow any update.

4 Task 2

The optimal policy is also displayed in figure 1. To run with the policy the following code snippet needs to be updated:

```
1 # TODO: Use the policy to take the optimal action (Task 2)
2 action = policy[state]
```

I rerun several episodes and observed the behavior. The boat is able to reach the harbour. However, it is not able to always reach the harbour. The sailor tends to take the narrow passage between the rocks. Due to the wind probability it sometimes crashes against the rocks.

5 Question 3

The sailor chose most of the times path through the narrow passage between the rocks. However, if we change the rock penalty to -10,

```
1 # Set up the environment
2 env = SailingGridworld(rock_penalty=-10)
```

the sailor tends to avoid narrow passage and goes most of the time below it.

6 Question 4

This question is conducted with rock penalty = -2. I rerun the algorithm with 10, 20, 30, 40, 50 iterations. For 10, 20, 30 iterations it didn't converge. However, for 40 iterations and higher it converged. The policy converged for 20 but not for 30 iterations. The policy needs less iterations to converge. That makes sense because as soon as it is figured which action take, the state values get fine tuned within the next iterations while the policy stays the same.

7 Task 3

The termination condition after convergence is integrated in the code:

```
1 max_diff_val = np.max(abs(abs(value_est) - abs(value_est_copy))
  )
2 if (max_diff_val < epsilon):
3     print("Converged! Value state converged in iteration: ", i
4         +1)
5     break
```

So as soon as the maximal updated difference is smaller than epsilon the value iteration stops.

8 Task 4

The code got updated to run the program for 1000 episodes and calculate the mean and variance of the discounted returns of the initial state. The mean is 0.6701 and the variance 1.3472.

```

1  # TODO: Run multiple episodes and compute the discounted
   returns (Task 4)
2  number_episodes = 1000
3  discounted_return_history = []
4
5  for i in range(number_episodes):
6      done = False
7      counter_discount = 0
8      discounted_return = 0
9
10     while not done:
11         # Select a random action
12         # TODO: Use the policy to take the optimal action (Task
           2)
13         action = policy[state]
14
15         # Step the environment
16         state, reward, done, _ = env.step(action)
17
18         # calculate discounted reward
19         discounted_return += reward * gamma**counter_discount
20         counter_discount += 1
21
22         # Render and sleep
23         #env.render()
24         #sleep(0.5)
25
26         discounted_return_history.append(discounted_return)
27         state = env.reset()
28
29  print("discounted return (initial state) - mean: ", np.mean(
       discounted_return_history))
30  print("discounted return (initial state) - std: ", np.std(
       discounted_return_history))

```

9 Question 5

In the book [1] we can see in equation (3.8) that the discounted returns depend on the rewards in each step and gamma. Furthermore, the state value depends on which policy we execute. The equation (3.12) states the value function of a state s under a policy π is the expected return when starting in s and following π thereafter. That's the relationship between the discounted return and the value function.

10 Question 6

No, the value iteration approach can't be applied to an unknown environment. It is a prerequisite to have complete knowledge of the environment which is stated in chapter 4.4 in book [1]. That includes probability of all transactions, actions, reward, and states. Then we can predict the next state value and reward. In this exercise the whole environment is known and therefore we can apply value iteration.

To be able to learn in an unknown environment is addressed with Monte Carlo estimation. The robot can learn with given experience from episodes.

References

- [1] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.