

Exercise 2

Christian Montecchiani - 100655834
ELEC-E8125 - Reinforcement Learning

September 26, 2022

1 Task 1

In the first task of this second assignment we were required to implement the value iteration algorithm. Below is reported my implementation.

```
1      ##### Your code starts here #####
2      temp_table = np.zeros((env.w, env.h))
3      for x in range(env.w):
4          for y in range(env.h):
5              value_list = np.zeros(env.n_actions)
6              for action in range(env.n_actions):
7                  transitions = env.transitions[x, y, action]
8                  action_value = 0
9                  for transition in transitions:
10                     next_state, reward, done, prob =
11                         transition.state, transition.reward,
12                         transition.done, transition.prob
13                     if done:
14                         action_value += prob * reward
15                         continue
16                     action_value += prob * (reward + gamma
17                         * v_est[next_state[0], next_state
18                             [1]])
19                     value_list[action] = action_value
20
21                 next_value = np.max(value_list)
22                 temp_table[x, y] = next_value
23                 policy[x, y] = np.argmax(value_list)
24
25      v_est = temp_table
26      env.draw_values_policy(v_est, policy)
27      ##### Your code ends here #####
```

1.1 Question 1.1

In this exercise the **agent** is represented by the *boat* and the **environment** is represented by the *grid world*.

1.2 Question 1.2

The state values of the harbor and rock states are zero. We obtain this result because they are *terminal states* and during the running of the algorithm they are never updated.

1.3 Question 1.3

If the penalty for hitting the rocks is set to -2 the sailor chooses the dangerous path between the rocks. If the reward for hitting the rocks is set to -10 then the sailor prefers to choose the longer but safer path.

2 Task 2

The state value function converges with a threshold set equal to 10^{-4} during the 42th iteration. This means that if I run the algorithm for 30 iterations the state value function still did not converge. Instead, if we check the policy, it converges faster. I check the convergence of the policy in two ways:

1. Check the changes between two consequent iterations and the policy converges at the 30th iteration.
2. Check the last time that the policy changes through all the iterations. It converges at iteration 32.

Generally, the policy needs less iteration to converge, which means that to find an optimal policy π^* is not required to compute the exact full value function.

```

1      ##### Your code starts here #####
2      temp_value = np.zeros((env.w, env.h))
3      temp_policy = np.zeros((env.w, env.h))
4
5      for x in range(env.w):
6          for y in range(env.h):
7              value_list = np.zeros(env.n_actions)
8
9              for action in range(env.n_actions):
10                 transitions = env.transitions[x, y, action]
11                 action_value = 0
12                 for transition in transitions:
13                     next_state, reward, done, prob =
14                         transition.state, transition.reward,
15                         transition.done, transition.prob
16                     if done:
17                         action_value += prob * reward
18                         continue
19                     action_value += prob * (reward + gamma
20                         * v_est[next_state[0], next_state
21                             [1]])
22                 value_list[action] = action_value
23
24                 next_value = np.max(value_list)
25                 temp_value[x, y] = next_value
26                 temp_policy[x, y] = np.argmax(value_list)
27
28     # Value function convergence
29     delta = np.abs(v_est - temp_value).max()
30     if delta < eps and i > 0 and first:
31         first = False
32         print(f"Value function has converged during
33             iteration: {i}")
34
35     v_est = temp_value
36
37     # Policy convergence
38     if np.array_equal(temp_policy, policy) == False:
39         last_change = i
40         policy = temp_policy
41
42     env.draw_values_policy(v_est, policy)
43
44     if i + 1 == iterations:
45         print(f"The last change in the policy is: {
46             last_change}")
47     ##### Your code ends here #####

```

3 Task 3

The number of iterations required for the value function to converge, when reward for crashing into the rocks is -2 , is 38. The code used is the same as the previous task.

4 Task 4

In this task we had to evaluate the learned policy and to do that we compute the *discounted return*:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

(as reported in [1]) of the initial state. The code is shown below.

```
1  # Eval policy
2  N = 1000 # TODO: change for task 4
3  returns = list()
4  for ep in range(N):
5      if (ep + 1) % 200 == 0:
6          print(f"Episodes {ep+1}/{N}")
7          state = env.reset()
8          done = False
9          rewards = list()
10         while not done:
11             ##### You code starts here #####
12
13             action = policy[state[0], state[1]]
14             # Take a step in the environment
15             state, reward, done, _ = env.step(action)
16             # Calculate discounted return for the initial state
17             rewards.append(reward)
18             if done:
19                 G = np.sum(np.array([(gamma**i) * rew for i,rew in
20                                     enumerate(rewards)]))
21                 returns.append(G)
22             ##### Your code ends here #####
```

The mean return and the standard deviation over 1000 episodes are: 0.66 and 1.36, respectively.

4.1 Question 4.1

As stated in [1] in Eq.3.12 the value function is:

$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Basically, the value function in a given state, s , is the expected discounted return started in that state and following policy π .

4.2 Question 4.2

In a reinforcement learning problem involving a robot exploring an *unknown* environment the value iteration approach cannot be applied. That is because the value iteration approach requires the to be in a fully observable Markov Decision Process (MDP). And this is the *unrealistic* the assumption that does not permits to use of the value iteration approach in each type of RL problems, since the majority of the environments, are not totally observable.

References

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