

Exercise 2

Eugeniu Vezeteu - 886240
ELEC-E8125 - Reinforcement Learning

September 21, 2020

1 Question 1 - What is the agent and the environment in this setup?

The agent is the sailor (boat), whose purpose is to reach the harbour. The environment is the sea, rocks, narrow passage between the rocks, wind and harbour.

2 Task 1 - Value iteration

The values are updated from the harbour to all the cells in the grid, closer to the harbour the higher the cell value is, see 1

r=0.00 V=0.53	r=0.00 V=0.59	r=0.00 V=0.66	r=0.00 V=0.73	r=0.00 V=0.66	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=0.00 V=7.72	r=0.00 V=8.71	r=0.00 V=9.84	r=10.00 V=0.00
r=0.00 V=0.59	r=0.00 V=0.66	r=0.00 V=0.74	r=0.00 V=0.83	r=0.00 V=0.80	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=0.00 V=6.95	r=0.00 V=7.80	r=0.00 V=8.76	r=0.00 V=9.84
r=0.00 V=0.65	r=0.00 V=0.73	r=0.00 V=0.83	r=0.00 V=0.94	r=0.00 V=1.06	r=0.00 V=1.21	r=0.00 V=1.74	r=0.00 V=2.39	r=0.00 V=3.21	r=0.00 V=4.20	r=0.00 V=5.59	r=0.00 V=6.19	r=0.00 V=6.95	r=0.00 V=7.80	r=0.00 V=8.71
r=0.00 V=0.64	r=0.00 V=0.72	r=0.00 V=0.81	r=0.00 V=0.90	r=0.00 V=1.01	r=0.00 V=1.14	r=0.00 V=1.64	r=0.00 V=2.24	r=0.00 V=2.97	r=0.00 V=3.84	r=0.00 V=4.98	r=0.00 V=5.52	r=0.00 V=6.19	r=0.00 V=6.95	r=0.00 V=7.72
r=0.00 V=0.65	r=0.00 V=0.73	r=0.00 V=0.81	r=0.00 V=0.90	r=0.00 V=0.82	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=0.00 V=4.91	r=0.00 V=5.51	r=0.00 V=6.18	r=0.00 V=6.84
r=0.00 V=0.73	r=0.00 V=0.82	r=0.00 V=0.92	r=0.00 V=1.02	r=0.00 V=1.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=-2.00 V=0.00	r=0.00 V=4.37	r=0.00 V=4.91	r=0.00 V=5.50	r=0.00 V=6.06
r=0.00 V=0.81	r=0.00 V=0.91	r=0.00 V=1.03	r=0.00 V=1.16	r=0.00 V=1.30	r=0.00 V=1.47	r=0.00 V=1.66	r=0.00 V=1.88	r=0.00 V=2.14	r=0.00 V=2.61	r=0.00 V=3.19	r=0.00 V=3.89	r=0.00 V=4.37	r=0.00 V=4.89	r=0.00 V=5.37
r=0.00 V=0.89	r=0.00 V=1.00	r=0.00 V=1.13	r=0.00 V=1.28	r=0.00 V=1.45	r=0.00 V=1.65	r=0.00 V=1.87	r=0.00 V=2.11	r=0.00 V=2.40	r=0.00 V=2.72	r=0.00 V=3.08	r=0.00 V=3.47	r=0.00 V=3.89	r=0.00 V=4.35	r=0.00 V=4.76
r=0.00 V=0.86	r=0.00 V=0.96	r=0.00 V=1.08	r=0.00 V=1.22	r=0.00 V=1.37	r=0.00 V=1.54	r=0.00 V=1.73	r=0.00 V=1.94	r=0.00 V=2.18	r=0.00 V=2.45	r=0.00 V=2.75	r=0.00 V=3.09	r=0.00 V=3.47	r=0.00 V=3.87	r=0.00 V=4.22
r=0.00 V=0.81	r=0.00 V=0.90	r=0.00 V=1.01	r=0.00 V=1.13	r=0.00 V=1.27	r=0.00 V=1.42	r=0.00 V=1.58	r=0.00 V=1.77	r=0.00 V=1.98	r=0.00 V=2.21	r=0.00 V=2.47	r=0.00 V=2.76	r=0.00 V=3.09	r=0.00 V=3.44	r=0.00 V=3.74

Figure 1: Task 1 - Value iteration for the sailor example after 100 iterations

$$V_{k+1}(s) = \max_a E[R_{t+1} + \gamma * V_k(s_{t+1} | s_t = s, a_t = a)] = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma * V_k(s')] \text{ from [1].}$$

```

1
2 def Task1(value_est,policy):
3     for iter in range(100):
4         env.clear_text()
5         #get all states
6         for i in range(env.w):
7             for j in range(env.h):
8                 history = []
9                 #compute the sum for all s'
10                for tran in env.transitions[i,j]:
11                    sum = 0.
12                    for s_,r,done,p in tran:
13                        sum += p * (r + (gamma * value_est[s_
14                            [0],s_[1]] if not done else 0))
15
16                    history.append(sum)
17
18                #update V(s), and pi(s)
19                value_est[i,j] = np.max(history)
20                policy[i,j] = np.argmax(history)
21
22    return value_est,policy

```

3 Question 2 - What is the state value of the harbour and rock states? Why?

The state value of the harbour and the rock states is 0. The value of the state is expected sum of future rewards, starting from that state and following policy π , that's why for harbour state value is zero (there are no more rewards to get) and for rock states (we never start from this state).

4 Task 2 - optimal policy

The sailor reaches the harbour most of the times, see Figure 2 On this task policy computed in the previous task was used.

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=-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
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=-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
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	r=0.00 V=8.71 a: Up
r=0.00 V=0.64 a: Right	r=0.00 V=0.72 a: Right	r=0.00 V=0.81 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	r=0.00 V=7.72 a: Up
r=0.00 V=0.65 a: Down	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=-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
r=0.00 V=0.73 a: Down	r=0.00 V=0.82 a: Down	r=0.00 V=0.92 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=-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
r=0.00 V=0.81 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: Right	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	r=0.00 V=5.37 a: Up
r=0.00 V=0.89 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	r=0.00 V=4.76 a: Up
r=0.00 V=0.86 a: Right	r=0.00 V=0.96 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	r=0.00 V=4.22 a: Up
r=0.00 V=0.81 a: Right	r=0.00 V=0.90 a: Right	r=0.00 V=1.01 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: Up	r=0.00 V=3.09 a: Up	r=0.00 V=3.44 a: Up	r=0.00 V=3.74 a: Up

Figure 2: Task 2 - Value iteration, Action and reward

```

1   for episod in range(10):
2       done = False
3       while not done:
4           # TODO: Use the policy to take the optimal action (
              Task 2)
5           action = policy[state]
6
7           # Step the environment
8           state, reward, done, _ = env.step(action)
9
10          # Render and sleep
11          env.render()
12          sleep(0.1)

```

5 Question 3 - Which path did the sailor choose?

With penalty of -2, the sailor take the dangerous path between the rocks (with is the short-est), but when changing the penalty to -10, the agent took the safe path below the rocks. We can notice from Figure 3 that the state values between the rocks are much smaller then in Figure 2.

r=0.00 V=0.41 a: Down	r=0.00 V=0.45 a: Down	r=0.00 V=0.50 a: Down	r=0.00 V=0.54 a: Down	r=0.00 V=0.48 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
r=0.00 V=0.46 a: Right	r=0.00 V=0.51 a: Down	r=0.00 V=0.57 a: Down	r=0.00 V=0.61 a: Down	r=0.00 V=0.54 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
r=0.00 V=0.52 a: Down	r=0.00 V=0.57 a: Down	r=0.00 V=0.64 a: Down	r=0.00 V=0.69 a: Down	r=0.00 V=0.63 a: Down	r=0.00 V=0.56 a: Left	r=0.00 V=-0.64 a: Left	r=0.00 V=0.40 a: Right	r=0.00 V=1.72 a: Right	r=0.00 V=3.32 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	r=0.00 V=8.71 a: Up
r=0.00 V=0.58 a: Down	r=0.00 V=0.64 a: Down	r=0.00 V=0.72 a: Down	r=0.00 V=0.78 a: Down	r=0.00 V=0.71 a: Down	r=0.00 V=0.64 a: Left	r=0.00 V=-0.60 a: Left	r=0.00 V=0.24 a: Right	r=0.00 V=1.48 a: Right	r=0.00 V=2.96 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	r=0.00 V=7.72 a: Up
r=0.00 V=0.65 a: Down	r=0.00 V=0.72 a: Down	r=0.00 V=0.81 a: Down	r=0.00 V=0.89 a: Down	r=0.00 V=0.79 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
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=1.01 a: Down	r=0.00 V=0.91 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
r=0.00 V=0.81 a: Right	r=0.00 V=0.91 a: Right	r=0.00 V=1.02 a: Down	r=0.00 V=1.14 a: Down	r=0.00 V=1.29 a: Down	r=0.00 V=1.45 a: Down	r=0.00 V=1.65 a: Down	r=0.00 V=1.87 a: Down	r=0.00 V=2.12 a: Down	r=0.00 V=2.40 a: Right	r=0.00 V=2.79 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	r=0.00 V=5.37 a: Up
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=1.44 a: Right	r=0.00 V=1.63 a: Right	r=0.00 V=1.85 a: Right	r=0.00 V=2.10 a: Right	r=0.00 V=2.38 a: Right	r=0.00 V=2.70 a: Right	r=0.00 V=3.06 a: Right	r=0.00 V=3.47 a: Right	r=0.00 V=3.89 a: Up	r=0.00 V=4.35 a: Up	r=0.00 V=4.76 a: Up
r=0.00 V=0.85 a: Right	r=0.00 V=0.96 a: Right	r=0.00 V=1.08 a: Right	r=0.00 V=1.21 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	r=0.00 V=4.22 a: Up
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=1.13 a: Right	r=0.00 V=1.26 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	r=0.00 V=3.74 a: Up

Figure 3: Question 3 - Value iteration, Action and reward with penalty of -10 for rocks

6 Question 4 - What happens if you run the algorithm for a smaller amount of iterations

With the smaller number of iterations the agent still can find the way to harbour, however this works with some number close to 100 iterations. Figure 4 shows that the agent cannot find the target with too small number of iterations (10 in this case).

As the number of iterations go to infinity, the value function and policy converge to optimal V^* and π^* .

Both, value function and policy need the same amount of iterations to converge, however, I will say that policy requires a bit more iterations than value function.

From the following formula we compute value function as max of expected reward for all actions. $V_{k+1}(s) = \max_a E[R_{t+1} + \gamma * V_k(s_{t+1}|s_t = s, a_t = a)]$ however, when we compute the policy we take the index from the maximum of the expected future rewards.

$\pi(s) = \operatorname{argmax}_a E[R_{t+1} + \gamma * V_k(s_{t+1}|s_t = s, a_t = a)]$, now imagine the case when there are more than one state with the same value function, in that case we will have more than one index of the max, and argmax will take one of them (wich may be the wrong), thats why policy need a bit more iterations.

r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 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
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=-1.10 a: Left	r=0.00 V=-0.08 a: Right	r=0.00 V=1.57 a: Right	r=0.00 V=3.30 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	r=0.00 V=8.71 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=-1.10 a: Left	r=0.00 V=-0.47 a: Right	r=0.00 V=1.19 a: Right	r=0.00 V=2.87 a: Right	r=0.00 V=4.96 a: Right	r=0.00 V=5.51 a: Up	r=0.00 V=6.19 a: Up	r=0.00 V=6.95 a: Up	r=0.00 V=7.72 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=0.00 V=4.90 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
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=-10.00 V=0.00 a: Left	r=0.00 V=4.31 a: Right	r=0.00 V=4.90 a: Up	r=0.00 V=5.50 a: Up	r=0.00 V=6.06 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Down	r=0.00 V=0.00 a: Down	r=0.00 V=0.00 a: Down	r=0.00 V=0.00 a: Down	r=0.00 V=0.00 a: Down	r=0.00 V=1.57 a: Right	r=0.00 V=3.57 a: Right	r=0.00 V=4.31 a: Up	r=0.00 V=4.88 a: Up	r=0.00 V=5.36 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=2.19 a: Up	r=0.00 V=3.57 a: Up	r=0.00 V=4.28 a: Up	r=0.00 V=4.73 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=2.17 a: Up	r=0.00 V=3.50 a: Up	r=0.00 V=4.06 a: Up	r=0.00 V=4.06 a: Up
r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=0.00 a: Left	r=0.00 V=2.03 a: Up	r=0.00 V=3.04 a: Up	r=0.00 V=3.04 a: Up

Figure 4: Question 4 - Value iteration, Action - Train for only 10 iterations

7 Task 3

```
1  #here is the updated function from the first task
2  def Task1(value_est,policy):
3      prev_value_est = value_est.copy()
4      for iter in range(100):
5          env.clear_text()
6          #get all states
7          for i in range(env.w):
8              for j in range(env.h):
9                  history = []
10                 #compute the sum for all s'
11                 for tran in env.transitions[i,j]:
12                     sum = 0.
13                     for s_,r,done,p in tran:
14                         sum += p * (r + (gamma * value_est[s_
15                             [0],s_[1]] if not done else 0))
16
17                     history.append(sum)
18
19                 #update V(s), and pi(s)
20                 value_est[i,j] = np.max(history)
21                 policy[i,j] = np.argmax(history)
22
23             if ((value_est - prev_value_est) < epsilon).all():
24                 print('Early stopping')
25                 break
26             else:
27                 prev_value_est = value_est.copy()
28
29     return value_est,policy
```

Result for task 3 is presented in Figure 5

8 Task 4

$$G_t = R_{t+1} + \gamma * R_{t+2} + \gamma^2 * R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k * R_{t+k+1}$$

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=-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
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=-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
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=-2.00 V=2.39 a: Right	r=-2.00 V=3.21 a: Right	r=-2.00 V=4.20 a: Right	r=-2.00 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	r=0.00 V=8.71 a: Up
r=0.00 V=0.64 a: Right	r=0.00 V=0.72 a: Right	r=0.00 V=0.81 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	r=0.00 V=7.72 a: Up
r=0.00 V=0.65 a: Down	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: Down	r=-2.00 V=0.00 a: Left	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	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
r=0.00 V=0.73 a: Down	r=0.00 V=0.82 a: Right	r=0.00 V=0.92 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: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	r=-2.00 V=0.00 a: Down	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
r=0.00 V=0.81 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: Right	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	r=0.00 V=5.37 a: Up
r=0.00 V=0.89 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	r=0.00 V=4.76 a: Up
r=0.00 V=0.86 a: Right	r=0.00 V=0.96 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	r=0.00 V=4.22 a: Up
r=0.00 V=0.81 a: Right	r=0.00 V=0.90 a: Right	r=0.00 V=1.01 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	r=0.00 V=3.74 a: Up

Figure 5: Task 3 - Run until convergence

```

1 history_G = []
2 for episod in range(1000):
3     state = env.reset()
4
5     done = False
6     G, k = 0, 0
7     while not done:
8         action = policy[state]
9         state, reward, done, _ = env.step(action)
10
11         G += (gamma**k)*reward
12         k+=1
13         #env.render()
14         #sleep(0.1)
15
16     history_G.append(G)
17 print('Avg:{}, Std:{}'.format(np.average(history_G), np.std(
    history_G)))

```

I got Avg:0.73, Std:1.35

9 Question 5 - What is the relationship between the discounted return and the value function?

$v_{\pi}(s) = E[G_t | s_t = s]$ Value function is the expected discounted return.

10 Question 6

If the environment is unknown, we cannot apply value iteration approach presented above. Because, we don't know the motion model $p(s', r|s, a)$, also, the state space may be infinite, or continuous. Using the value function in the real world scenario may be very computational expensive.

References

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