

# Reinforcement Learning - Exercise 2

Muhammad Kamal Memon  
600442

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

The agent in this environment is the sailor who wants to reach the harbour. The environment consists of the sea, the harbour, the rocks and the narrow passageway in between them. Also the heavy wind conditions and the calmness of the sea are part of the environment.

## Task 1

```
# TODO: Compute value function and policy.
value_est, policy = np.zeros((env.w, env.h)), np.zeros((env.w, env.h)),

def action_values(state, V):
    A = []
    for transition in env.transitions[state[0], state[1]]:
        action_value = 0
        for next_state, reward, done, prob in transition:
            action_value += prob * (reward + (gamma * V[next_state] if not done else 0))
        A.append(action_value)
    return A

for lp in range(100):
    env.clear_text()
    for w in range(env.w):
        for h in range(env.h):
            A = action_values((w, h), value_est)
            value_est[w][h] = np.max(A)
            policy[w][h] = np.argmax(A)
    env.draw_values(value_est)
    env.draw_actions(policy)
    env.render()
```

The values get updated from the harbour state all the way down to bottom left. Figure 1 below shows how the values are updated after running the value iteration algorithm:

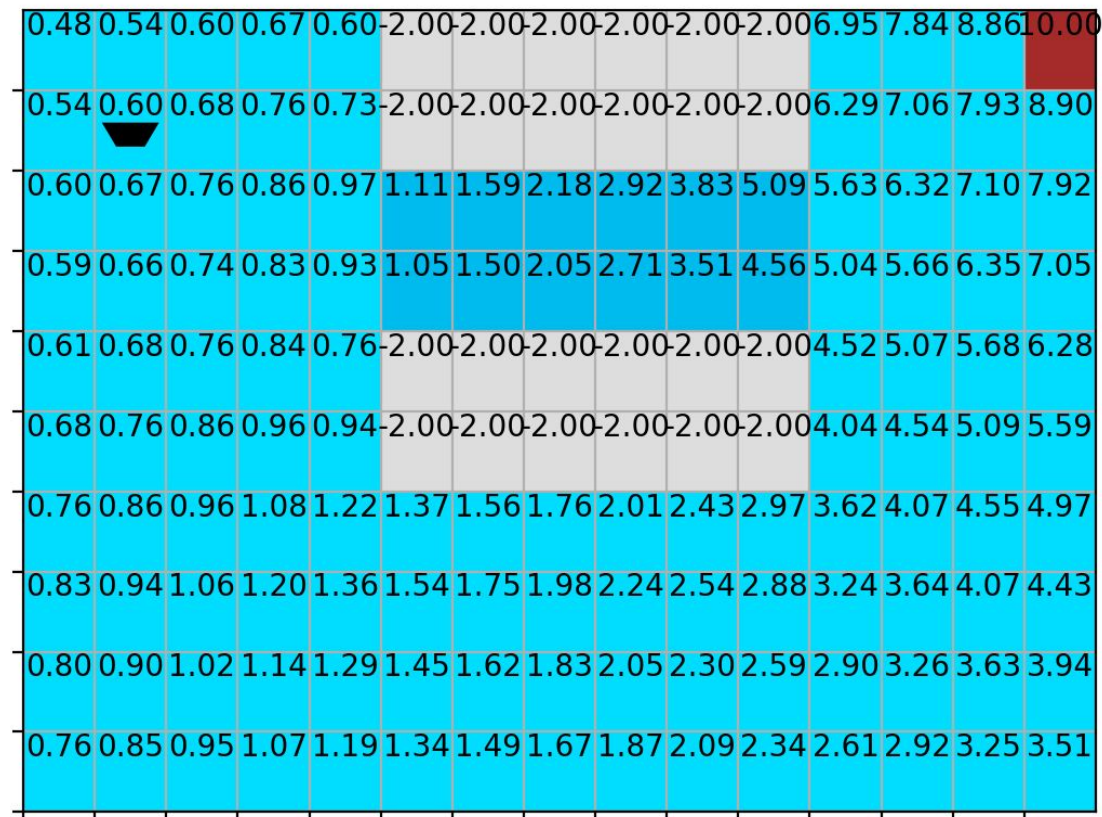


Figure1: Values in the grid state space after Value Iteration

## Task 2

```

sleep(3)
done = False
while not done:
    # TODO: Use the computed policy here.
    action = policy[state]
    state, reward, done, _ = env.step(action)
    env.render()
    sleep(0.5)

```

The policy computed above in Task1 is implemented and the sailor manages to reach the goal state most of times the algorithm was run.

0.48	0.54	0.60	0.67	0.60	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	6.95	7.84	8.86	10.00
Right	Right	Down	Down	Left	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Right	Stay
0.54	0.60	0.68	0.76	0.73	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	6.29	7.06	7.93	8.90
Right	Right	Down	Down	Down	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Up	Up
0.60	0.67	0.76	0.86	0.97	1.11	1.59	2.18	2.92	3.83	5.09	5.63	6.32	7.10	7.92	
Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up
0.59	0.66	0.74	0.83	0.93	1.05	1.50	2.05	2.71	3.51	4.56	5.04	5.66	6.35	7.05	
Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up	Up
0.61	0.68	0.76	0.84	0.76	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	4.52	5.07	5.68	6.28
Down	Right	Down	Down	Left	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Up	Up	Up	Up
0.68	0.76	0.86	0.96	0.94	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	-2.00	4.04	4.54	5.09	5.59
Down	Right	Down	Down	Down	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Up	Up	Up	Up
0.76	0.86	0.96	1.08	1.22	1.37	1.56	1.76	2.01	2.43	2.97	3.62	4.07	4.55	4.97	
Right	Right	Right	Right	Right	Down	Down	Down	Down	Right	Right	Right	Up	Up	Up	Up
0.83	0.94	1.06	1.20	1.36	1.54	1.75	1.98	2.24	2.54	2.88	3.24	3.64	4.07	4.43	
Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up	Up
0.80	0.90	1.02	1.14	1.29	1.45	1.62	1.83	2.05	2.30	2.59	2.90	3.26	3.63	3.94	
Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up	Up
0.76	0.85	0.95	1.07	1.19	1.34	1.49	1.67	1.87	2.09	2.34	2.61	2.92	3.25	3.51	
Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up

Figure 2: Value function, policy and sailor end

## Question 2

The sailor always chooses the path through the narrow passage when the reward to hit the rocks is -2. However when the reward is set to -10 it decides to not take the path from the narrow passage as now the risk of getting negative reward is high and so it goes the other way around to reach the goal.

## Question 3

We ran the algorithm for 30 and 15 iterations instead of 100. For 30 iterations the value function and policy still converges and the sailor is able to find its way to the harbour. On the other hand for 15 iterations they don't converge. In more experiments it resulted so that the policy gets converged but the value function take more iterations to do so. The reason being that the value function need the whole state space to take into account while policy gets converged in less iterations. See the Figure below for 15 iterations:



0.00	0.00	0.00	0.00	0.00	10.00	10.00	10.00	10.00	10.00	10.00	10.00	6.95	7.84	8.86	10.00
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Right	Stay
0.00	0.00	0.00	0.00	0.00	10.00	10.00	10.00	10.00	10.00	10.00	10.00	6.29	7.06	7.93	8.90
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.59	0.39	1.59	3.04	5.09	5.63	6.32	7.10	7.92
Stay	Stay	Stay	Stay	Stay	Left	Right	Right	Right	Right	Right	Right	Right	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	-0.68	0.25	1.37	2.72	4.56	5.04	5.66	6.35	7.05	
Stay	Stay	Stay	Stay	Stay	Left	Right	Right	Right	Right	Right	Right	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	10.00	10.00	10.00	10.00	10.00	10.00	10.00	4.52	5.07	5.68	6.28
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	10.00	10.00	10.00	10.00	10.00	10.00	10.00	4.04	4.54	5.09	5.59
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.78	1.62	2.16	2.61	3.62	4.07	4.55	4.97
Stay	Stay	Stay	Stay	Stay	Left	Down	Down	Down	Down	Down	Right	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.91	1.85	2.44	2.85	3.24	3.64	4.07	4.43
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Right	Right	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.04	1.97	2.51	2.89	3.26	3.63	3.94
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Right	Up	Up	Up	Up
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.05	1.99	2.53	2.90	3.24	3.50
Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Stay	Right	Right	Right	Up	Up	Up

Figure3: Incomplete value function and policy did not converge

### Task 3

```
def action_values(state, V):
    A = []
    for transition in env.transitions[state[0], state[1]]:
        action_value = 0
        for next_state, reward, done, prob in transition:
            action_value += prob * (reward + (gamma *
                                           V[next_state] if not done else 0))
        A.append(action_value)
    return A

old_values = value_est.copy()
eps = 0.0001

for lp in range(100):
    env.clear_text()
    for w in range(env.w):
        for h in range(env.h):
            A = action_values((w, h), value_est)
            value_est[w][h] = np.max(A)
            policy[w][h] = np.argmax(A)
    if ((value_est - old_values) < eps).all():
```

```

        print('Stopping')
        break
    else:
        old_values = np.copy(value_est)
        env.draw_values(value_est)
        env.draw_actions(policy)
        env.render()

```

## Task 4

```

discounted_reward = 0
i = 0
done = False
while not done:
    # TODO: Use the computed policy here.
    action = policy[state]
    state, reward, done, _ = env.step(action)
    discounted_reward += (gamma ** i) * reward
    i += 1
    env.render()
    sleep(0.5)
print(discounted_reward)

```

Discounted return	Sailor status
1.0941	Managed to reach the harbour
-0.4117	Hit the rocks
1.8530	Managed to reach the harbour
-0.9565	Hit the rocks
1.5009	Managed to reach the harbour
1.6677	Managed to reach the harbour

## Question 4

For any policy  $\pi$  and any state  $s$ , the following holds:

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

That is, the value function is expectation (following a particular policy) of the discounted return  $G_t$  for any state  $s$  [Chapter 3 Sutton & Barto]

### **Question 5**

The value iteration approach implemented here wouldn't fit too well to a real world scenario. Mainly because here we have finite state space with a deterministic set of actions. Moreover, in a real world setting computing the value function the way we did can also get very computationally expensive as we are maximizing over action values. Also the model of the problem would not be this straight forward.

### **References:**

The exercise involved discussion with following students:

- Gadidjah Ogmundsdottir
- Hector Laria Matecon



