# Exercise 2

Nguyen Xuan Binh - 887799
ELEC-E8125 - Reinforcement Learning

July 16, 2023

## Task 1 - 30 points

**i) Implement value iteration for the sailor example in file value_iteration.ipynb, assuming the discount factor value $\gamma = 0.9$.**

**ii) In addition to the state values, compute the policy – path to the harbour, using computed state values. Run your implementation for 100 iterations. Attach this .pkl file with the estimated state values and policy into your submission.**

This is the converged value and policy after 100 iterations and the Python code

```python
def get_values_policy(iterations):
    v_est = np.zeros((env.w, env.h))
    policy = np.zeros((env.w, env.h))
    # env.draw_values_policy(v_est, policy)

    value_converge_index = -1
    policy_converge_index = -1
    check_value_converge = False
    check_policy_converge = False

    for i in range(1, iterations + 1):
        # TODO: Task 1, implement the value iteration and
            policy
        # TODO: Task 2, convergency of the value function and
            policy
        old_v_est = v_est.copy()
        new_v_est = v_est.copy()
        old_policy = policy.copy()
        new_policy = np.zeros((env.w, env.h))

        # Iterate over all state coordinates
        for x in range(env.w):
            for y in range(env.h):
                Q_values = np.zeros(env.n_actions)
                # Iterate over all actions and calculate the Q
                    value for each action
                for action in range(env.n_actions):
                    transitions = env.transitions[x, y, action]
                    Q_value = 0
                    for next_state, reward, done, probability
                        in transitions:
                        if not done:
                            next_state_value = v_est[next_state
                                ]
                        else:
                            next_state_value = 0 # Terminal
                                state value should be zero
                        Q_value += probability * (reward +
                            gamma * next_state_value)
                    Q_values[action] = Q_value
                new_v_est[x, y] = np.max(Q_values) # Updating
                    state value with maximum Q value
                new_policy[x, y] = np.argmax(Q_values) #
                    Updating state policy with action that
                    maximizes Q values
```

```
1       # Calculate  the  change  in  the  value  function  among  all
           states
2         deltas = np.abs(new_v_est - old_v_est)
3
4         # Check  if  the  value  function  has  converged
5         if np.all(deltas < eps):
6             if not check_value_converge:
7                 value_converge_index = i
8                 check_value_converge = True
9             # The  draw_values_policy  function  is  slow  and
                   called  only  once
10            # env.draw_values_policy(new_v_est,  new_policy)
11
12        # Check  if  the  policy  function  has  converged
13        if np.array_equal(old_policy, new_policy):
14            if not check_policy_converge:
15                policy_converge_index = i
16                check_policy_converge = True
17
18        #print(f"Iteration  {i+1}  finishes")
19
20        # Updating  the  value  function  and  policy
21        v_est = new_v_est
22        policy = new_policy
23        #print(f"Iteration  {i}  finishes")
24
25        # The  draw_values_policy  function  is  slow  and  called
                   only  once
26        if i == iterations:
27            env.draw_values_policy(v_est, policy)
28        #print(i)
29        #print(iterations)
30
31    print(f"Value iteration has converged after {
          value_converge_index} iterations")
32    print(f"Policy function has converged after {
          policy_converge_index} iterations")
33
34    return v_est, policy
```

## Question 1.1 - 5 points

**What is the agent and the environment in this sailor gridworld?**

The agent is the sailor boat and the environment is the whole gridworld, consisting of the sea, the harbour and the rocks.

# Question 1.2 - 5 points

**What is the state value of the harbour and rock states? Why?**

The harbour and the rocks are the terminal states. Therefore, their values should not contribute to updating the values of the sea states, so the values of the harbour and rock states are zero.

# Question 1.3 - 5 points

**Which path did the sailor choose, the safe path below the rocks, or the dangerous path between the rocks? If you change the reward for hitting the rocks to -10 (that is, make the sailor value life more), does he still choose the same path?**

The path the sailor chose is the dangerous path between the rocks when the terminal state value of the rocks is -2. However, if we change the reward (or penalty) for hitting the rocks to -10, the sailor no longer chose the dangerous path. Instead, he chose the safer detour around the rocks. This is probably because the reward of reaching the harbour is not high enough to mitigate the risk of travelling through the dangerous passage between the rocks.

# Task 2 - 15 points

**What happens if you run the algorithm for 30 iterations? Do the value function and policy still converge? For the value function, you can assume they have converged if the maximum change in value is lower than a certain threshold $\epsilon = 10^{-4}$ :**

$$max_s |V_k(s) - V_{k-1}(s)| < \epsilon, (2)$$

**where $V_k(s)$ is the estimated value of state s in k-th iteration of the algorithm. Generally, which of them - the policy or value function - needs less iterations to converge, if any? Justify your answer.**

When I run the algorithms for 30 iterations, the value estimation has not managed to converge within the allowed threshold, while the policy still manages to converge after 26 iterations. Generally, the policy function needs less iterations to converge than the value function. This is because the policy space of the sailor world is relatively small (4 actions), so the algorithm can find a clear strategy structure quickly, while the value iteration is a continuous float number, which takes more iterations to get to the correct decimal precision.

# Task 3 - 5 points

**Set the reward for crashing into the rocks back to -2. Change the termination condition of your algorithm to make it run until convergence, see Eq. (2). Report the number of iterations required for the value function to converge.**

The number of iterations required for the value function to converge is 39 (my iterations start counting from 1, not 0).

# Task 4 - 10 points

**Evaluate your learned policy for N = 1000 episodes, and compute the discounted return of the initial state, see [1] Eq. (3.8), for each episode. The reward for crashing into rocks must be kept at -2 for this exercise. Report the average and standard deviation of the initial state's discounted return over the N=1000 episodes.**

The average of the initial state's discounted return discounted return is 0.6437
The standard deviation of the initial state's discounted return discounted return is 1.3685

# Question 4.1 - 10 points

**What is the relationship between the discounted return and the value function? Explain briefly.**

Definition: The value function is the expected cumulative reward of each state when the agent follows the policy starting from the initial state. Additionally, the discounted return is the sum of the rewards accumulated over N steps, with each reward discounted by a factor $\gamma$ after each step. Therefore, the value function at a state s is the expected sum of discounted rewards if the agent starts from the known initial state and follows a fixed policy.

# Question 4.2 - 15 points

**Imagine a reinforcement learning problem involving a robot exploring an unknown environment. Could the value iteration approach used in this exercise be applied directly to that problem? Why/why not? Which assumptions are unrealistic, if any?**

My answer is yes and no. Yes in that the principles and the environment settings are the same, but no because some assumptions in this exercise is too simple for the RL problem involving a robot exploring unknown environments. The unrealistic assumptions in this exercise that cannot be applied to the robot problem is:
- Four discrete actions: Left, Right, Up, Down. In reality, the robot movement is much more flexible, which could be described as a continuous domain instead of the discrete domain.
- Value iteration algorithm assumes that the agent has full access to the state space and and transition probabilities. In reality, these information could not be easily be observed.
- the searching space would be very large (which could have more than million states), rendering the value iteration too slow and impractical.