

## ✓ CS640 Homework 5: MDP, RL and Adversarial Search

This assignment covers three topics shown in the title. There is a problem for each topic.

In particular, the MDP problem asks you to implement both the value iteration and the policy iteration, the RL problem asks you to implement the Q-learning algorithm, and in the adversarial search problem, you need to run the min-max algorithm by hand.

### Collaboration

You must answer all questions **independently**, including the coding ones.

### Instructions

#### General Instructions

In an ipython notebook, to run code in a cell or to render [Markdown+LaTeX](#) press `Ctrl+Enter` or `[>|]` (like "play") button above. To edit any code or text cell (double) click on its content. To change cell type, choose "Markdown" or "Code" in the drop-down menu above.

Most of the written questions are followed up a cell for you enter your answers. Please enter your answers in a new line below the **Answer** mark. If you do not see such cell, please insert one by yourself. Your answers and the questions should **not** be in the same cell.

#### Instructions on Math

Some questions require you to enter math expressions. To enter your solutions, put down your derivations into the corresponding cells below using LaTeX. Show all steps when proving statements. If you are not familiar with LaTeX, you should look at some tutorials and at the examples listed below between  $$.$. The [OEIS website](#) can also be helpful.$

Alternatively, you can scan your work from paper and insert the image(s) in a text cell.

### Submission

Once you are ready, save the note book as PDF file (File -> Print -> Save as PDF) and submit via Gradescope.

Please select pages to match the questions on Gradescope. **You may be subject to a 5% penalty if you do not do so.**

✓ Q1: MDP

➤ Q1.0: Setups

[ ] ↳ 4 cells hidden

✓ Q1.1: Value iteration

Implement value iteration by finishing the following function, and then run the cell.

```

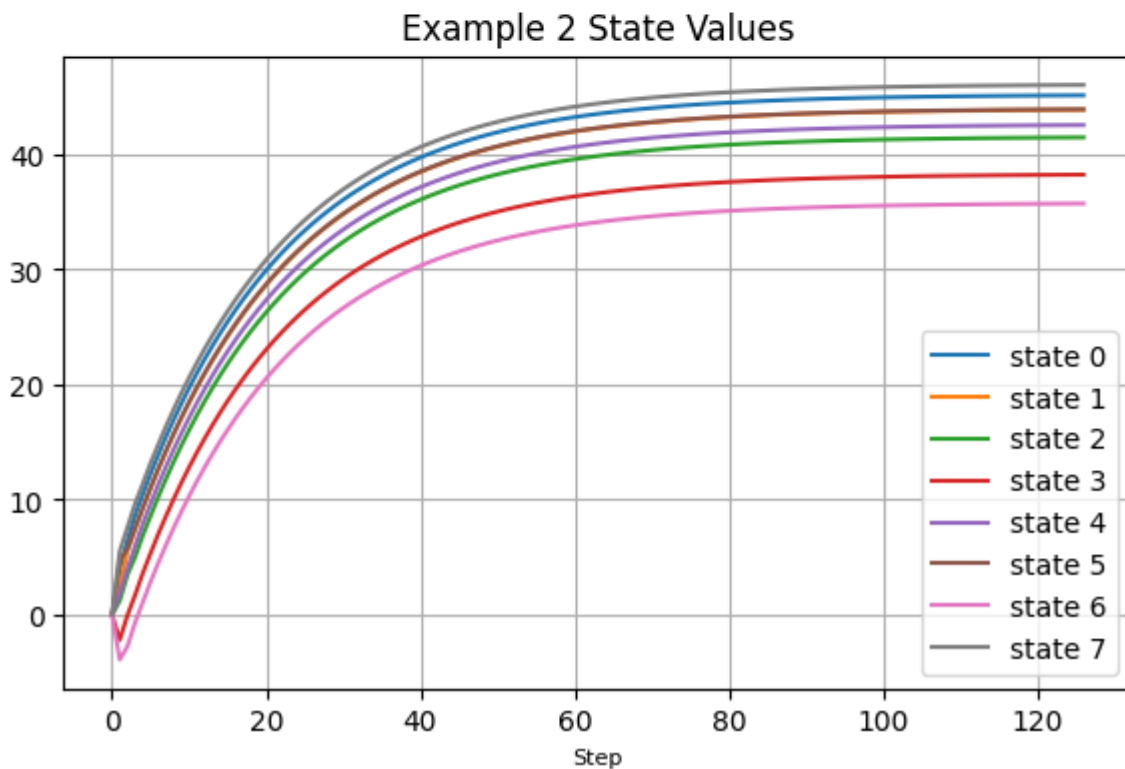
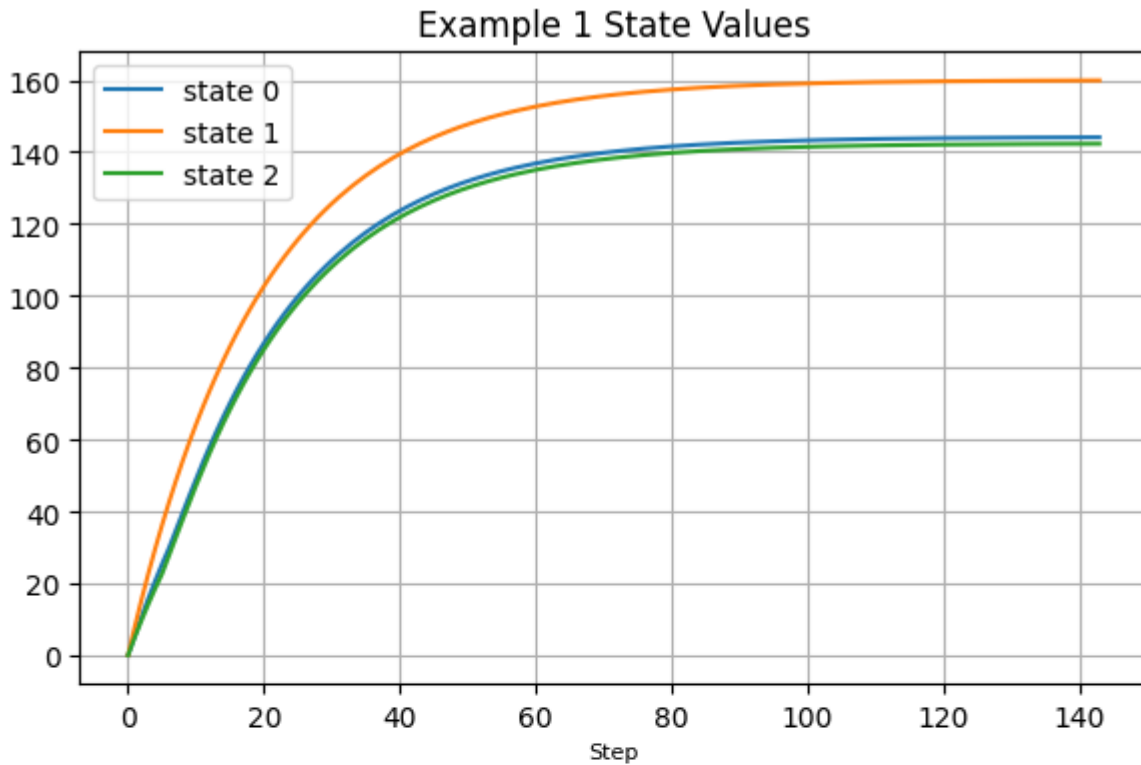
1 def value_iteration(states, actions, T, R, gamma = 0.95, tolerance = 1e-2, max_steps = 1000)
2     Vs = [] # all state values
3     Vs.append(np.zeros(len(states))) # initial state values
4     steps, convergent = 0, False
5     while not convergent and steps < max_steps: # complete this loop
6         ##### start of your code #####
7         # compute new state values as an array, and append the array to the list Vs
8         V = -np.ones(len(states)) * sys.maxsize
9         for s in states:
10            for a in actions:
11                new_value = np.sum(T[s, a] * (R[s, a] + gamma * Vs[-1]))
12                V[s] = max(V[s], new_value)
13            Vs.append(V)
14
15            ##### end of your code #####
16            convergent = np.linalg.norm(Vs[-1] - Vs[-2]) < tolerance
17            steps += 1
18        if steps == max_steps:
19            print("Max iterations reached before convergence.")
20            sys.exit(1)
21        ##### start of your code #####
22        # compute Q from the last V array
23        # extract the optimal policy and name it "policy" to return
24        Q = np.zeros((len(states), len(actions)))
25        for s in states:
26            for a in actions:
27                Q[s, a] = np.sum(T[s, a] * (R[s, a] + gamma * Vs[-1]))
28        policy = np.argmax(Q, axis = 1)
29        ##### end of your code #####
30        return np.array(Vs), policy, steps
31
32 states, actions, T, R = example_1
33 gamma, tolerance, max_steps = 0.95, 1e-2, 1000
34 Vs_1, policy_1, steps_1 = value_iteration(states, actions, T, R, gamma, tolerance, max_steps)
35
36 states, actions, T, R = example_2
37 gamma, tolerance, max_steps = 0.95, 1e-2, 1000
38 Vs_2, policy_2, steps_2 = value_iteration(states, actions, T, R, gamma, tolerance, max_steps)
39
40 fig, axes = plt.subplots(2, 1, figsize = (6, 8))
41 for i in range(Vs_1.shape[1]):
42     axes[0].plot(Vs_1[:, i], label = "state " + str(i))
43 axes[0].set_title("Example 1 State Values", fontsize = 12)
44 axes[0].set_xlabel("Step", fontsize = 8)
45 axes[0].legend()
46 axes[0].grid(True)
47 for i in range(Vs_2.shape[1]):
48     axes[1].plot(Vs_2[:, i], label = "state " + str(i))
49 axes[1].set_title("Example 2 State Values", fontsize = 12)
50 axes[1].set_xlabel("Step", fontsize = 8)
51 axes[1].legend()

```

```

52 axes[1].grid(True)
53 fig.tight_layout()
54 plt.show()
55
56 print("Optimal policy for example 1: " + str(policy_1))
57 print("Optimal policy for example 2: " + str(policy_2))

```



Optimal policy for example 1: [1 0 1]  
 Optimal policy for example 2: [3 4 4 3 1 1 4 1]

## ✓ Q1.2: Policy iteration

Implement policy iteration by finishing the following function, and then run the cell.

```

1 def policy_iteration(states, actions, T, R, gamma = 0.95, tolerance = 1e-2, max_steps = 1
2     policy_list = [] # all policies explored
3     initial_policy = np.array([np.random.choice(actions) for s in states]) # random polic
4     policy_list.append(initial_policy)
5     Vs = [] # all state values
6     Vs = [np.zeros(len(states))] # initial state values
7     steps, convergent = 0, False
8     while not convergent and steps < max_steps:
9         ##### start of your code #####
10        # 1. Evaluate the current policy, and append the state values to the list Vs
11        convergent_eval = False
12        V = Vs[-1]
13        while not convergent_eval:
14            V_new = np.zeros(len(states))
15            for s in states:
16                a = policy_list[-1][s]
17                Q = 0
18                for s_next in states:
19                    Q += T[s, a, s_next] * (R[s, a, s_next] + gamma * V[s_next])
20                V_new[s] = Q
21            convergent_eval = np.linalg.norm(V_new - V) < tolerance
22            V = V_new
23        Vs.append(V_new)
24
25        # 2. Extract the new policy, and append the new policy to the list policy_list
26        Q_values = -np.ones((len(states), len(actions))) * sys.maxsize
27        for s in states:
28            for a in actions:
29                Q_values[s, a] = np.sum(T[s, a] * (R[s, a] + gamma * Vs[-1]))
30        new_policy = np.argmax(Q_values, axis = 1)
31        policy_list.append(new_policy)
32        ##### end of your code #####
33        steps += 1
34        convergent = (policy_list[-1] == policy_list[-2]).all()
35    if steps == max_steps:
36        print("Max iterations reached before convergence.")
37        sys.exit(1)
38    return Vs, policy_list, steps
39
40 print("Example MDP 1")
41 states, actions, T, R = example_1
42 gamma, tolerance, max_steps = 0.95, 1e-2, 1000
43 Vs, policy_list, steps = policy_iteration(states, actions, T, R, gamma, tolerance, max_st
44 for i in range(steps):
45     print("Step " + str(i))
46     print("state values: " + str(Vs[i]))
47     print("policy: " + str(policy_list[i]))
48     print()
49 print()
50 print("Example MDP 2")
51 states, actions, T, R = example_2

```

```

52 gamma, tolerance, max_steps = 0.95, 1e-2, 1000
53 Vs, policy_list, steps = policy_iteration(states, actions, T, R, gamma, tolerance, max_st
54 for i in range(steps):
55     print("Step " + str(i))
56     print("state values: " + str(Vs[i]))
57     print("policy: " + str(policy_list[i]))
58     print()

```

Example MDP 1

Step 0

state values: [0. 0. 0.]

policy: [1 1 1]

Step 1

state values: [43.65839589 40.5921489 39.70489023]

policy: [1 0 0]

Step 2

state values: [112.95002513 159.85584651 99.92720955]

policy: [1 0 1]

Example MDP 2

Step 0

state values: [0. 0. 0. 0. 0. 0. 0. 0.]

policy: [2 4 2 0 3 3 2 0]

Step 1

state values: [ -7.51257276 -3.62936276 -9.23621336 -9.75080134 -8.63433687

-5.99069045 -14.67723454 -3.83794082]

policy: [3 4 4 3 1 1 4 1]

## > Q1.3: More tests

The following block tests both of your implementations with even more random MDPs. Simply run the cell.

[ ] ↴ 1 cell hidden

## ✓ Q2: RL

### > Q2.0: Setups

[ ] ↴ 13 cells hidden

## ✓ Q2.1: Q-learning

### Implement Q-learning

Implement Q-learning by modifying the specified part in the following block. Currently it is only taking a random action.

```
1 def Q_learning(env, episodes = 100000, alpha = 0.1, gamma = 0.6, epsilon = 0.1):
2     Q_values = np.zeros([env.observation_space.n, env.action_space.n])
3     rewards_list = [0] * episodes
4     for episode in range(episodes):
5         state = env.reset()
6         done = False
7         while not done:
8             #####
9             # TO DO: implement Q-learning update
10            if np.random.uniform(0, 1) < epsilon:
11                action = env.action_space.sample()
12            else:
13                action = np.argmax(Q_values[state])
14            next_state, reward, done, _, info = env.step(action) # you shouldn't change th
15            Q_values[state, action] = (1 - alpha) * Q_values[state, action] + alpha * (rew
16            ##### End of your code #####
17            rewards_list[episode] += reward
18            state = next_state
19     return Q_values, rewards_list
```

## ✓ Q3

### Q3.1

6, 6, 3, 6, 10, 3, 8, 20, 6, 10, 15, 7, 3, 16, 8

### Q3.2

0-1-3-8-18

### Q3.3

6, 13, 14, 20, 22, 27, 28, 29, 30