```
1
    #Nihal Ranchod - 2427378
 2
    #Lisa Godwin -2437980
 3
    import numpy as np
 4
 5
    import matplotlib.pyplot as plt
 6
    import seaborn as sns
 7
    class GridworldMDP:
 8
 9
         def init (self):
             self.grid_size = 4
10
             self.initial_state = (3, 0)
11
             self.goal state = (0, 0)
12
             self.actions = ['up', 'down', 'left', 'right']
13
14
             self.state = self.initial_state
15
         def step(self, state, action):
16
17
             if action not in self.actions:
                 raise ValueError("Invalid action")
18
19
20
             x, y = state
             if action = 'up':
21
                 x = \max(x - 1, 0)
22
             elif action = 'down':
23
                 x = min(x + 1, self.grid_size - 1)
24
25
             elif action = 'left':
                 y = \max(y - 1, 0)
26
             elif action = 'right':
27
                 y = min(y + 1, self.grid_size - 1)
28
29
             new_state = (x, y)
30
31
             reward = -1
32
             if new state = self.goal state:
33
                 reward = 0
34
35
             return new state, reward
36
37
         def is_terminal(self, state):
             return state = self.goal state
38
39
    # Policy Evaluation in Place:
40
    # The value function is updated immediately after evaluating each state.
41
42
    # The new value of a state is used in the subsequent evaluations of other
    states within the same iteration.
    def policy_evaluation_in_place(env, gamma, theta):
43
         V = np.zeros((env.grid size, env.grid size))
44
45
         iterations = 0
         while True:
46
             delta = 0
47
             for x in range(env.grid size):
48
49
                 for y in range(env.grid_size):
50
                     state = (x, y)
                     if env.is terminal(state):
51
                         continue
52
```

```
53
                      v = V[state]
 54
                      V[state] = sum(1/len(env.actions) * (reward + gamma *
     V[new state])
                                      for action in env.actions
 55
                                      for new_state, reward in [env.step(state,
 56
     action)])
                      delta = max(delta, abs(v - V[state]))
 57
              iterations += 1
 58
              if delta < theta:</pre>
 59
 60
                  break
          return V, iterations
 61
 62
 63
     # Policy Evaluation with Two Arrays:
 64
     # A temporary array (V_new) is used to store the updated values of the
     states.
     # The value function V is only updated after all states have been evaluated
 65
      in the current iteration
 66
      def policy_evaluation_two_array(env, gamma, theta):
          V = np.zeros((env.grid size, env.grid size))
 67
 68
          iterations = 0
          while True:
 69
              delta = 0
 70
 71
              V \text{ new = np.copy}(V)
              for x in range(env.grid_size):
 72
 73
                  for y in range(env.grid_size):
 74
                      state = (x, y)
                      if env.is_terminal(state):
 75
 76
                           continue
                      V new[state] = sum(1/len(env.actions) * (reward + gamma *
 77
     V[new state])
                                          for action in env.actions
 78
 79
                                          for new state, reward in [env.step(state,
     action)])
 80
                      delta = max(delta, abs(V_new[state] - V[state]))
 81
              V = V new
 82
              iterations += 1
              if delta < theta:</pre>
 83
 84
                  break
 85
          return V, iterations
 86
 87
     def main():
          env = GridworldMDP()
 88
 89
          theta = 0.01
          gammas = np.logspace(-0.2, 0, num=20)
 90
 91
 92
          iterations in place = []
          iterations_two_array = []
 93
 94
 95
          for gamma in gammas:
              _, it_in_place = policy_evaluation_in_place(env, gamma, theta)
 96
 97
              _, it_two_array = policy_evaluation_two_array(env, gamma, theta)
              iterations in place.append(it in place)
 98
              iterations two array.append(it two array)
 99
100
          V, = policy evaluation in place(env, 1, theta)
101
```

```
plt.figure(figsize=(6, 5))
102
103
          sns.heatmap(V, annot=True, cmap='crest', cbar=True)
          plt.title('Value Function Heatmap for $\gamma$ = 1')
104
          plt.show()
105
106
          plt.figure(figsize=(8, 6))
107
          plt.plot(gammas, iterations_in_place, label='Policy Evaluation In-place',
108
     color="slateblue")
         plt.plot(gammas, iterations_two_array, label='Policy Evaluation Two-
109
     array', color="deeppink")
         plt.xlabel('Discount Rate $\gamma$')
110
111
          plt.ylabel('Iterations to Convergence')
         plt.title('Convergence Rates for Policy Evaluation Algorithms')
112
          plt.legend()
113
         plt.show()
114
115
     if __name__ = '__main__':
116
117
         main()
```



