

§ MDP §

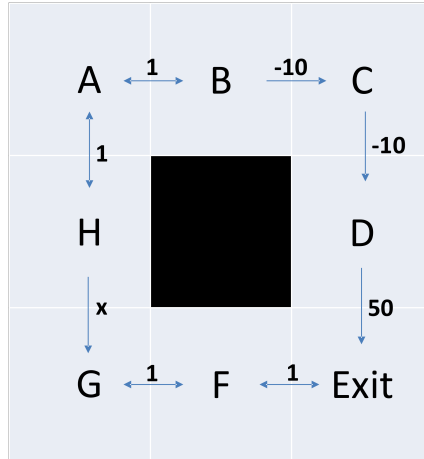


Figure 1: A 3x3 grid world with the center cell not accessible (colored in shade), there are states A through H, and Exit.

Given a MDP representing a 3x3 grid world with the center cell not accessible (colored in shade), there are states A through H, and Exit. The actions at each state include moving to its neighbors as indicated by the arrow and exit. The number along with the arrow is the immediate reward the agent would receive [Pay attention, the reward from H to G is x]. For exit action, there is no reward.

Codes are publicly available at [GitHub: MDP-Problem-AI-CSI-5130](#)

Problem 1: Value Iteration with $\gamma=1$

Assume all the actions are deterministic. If the discount factor $\gamma=1$, please calculate the $V^*(A)$, $V^*(C)$, $V^*(G)$, $V^*(Exit)$.

Answer:

I've written a code (Appendix I) for calculation of the value of each state and finding the optimal policy. When $\gamma=1$, the value iteration does not converge and the code output would be like this for 3 iterations:

Code output for Problem 1

```
1 After 3 iterations:
2 The final value for each state is:
3 {'A': 3, 'B': 30, 'C': 41, 'D': 52, 'E': 3, 'F': 3, 'G': 3, 'H': 3, '0': 0}
4 And, the optimal policy is:
5 {'A': 'r', 'B': 'r', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
  ↪ '0': 'r'}
```

Problem 2: Value Iteration with $\gamma=0.5$

Assume all the actions are deterministic. If the discount factor $\gamma=0.5$, please calculate the $V^*(A)$, $V^*(C)$, $V^*(G)$, $V^*(Exit)$.

Answer:

For $\gamma=0.5$, the value iteration will be converged in 55 iterations and the code output would be like this:

Code output for Problem 2

```

1 After 55 iterations:
2 The final value for each state is:
3 {'A': 2.0, 'B': 2.0, 'C': 15.5, 'D': 51.0, 'E': 2.0, 'F': 2.0, 'G': 2.0, 'H': 2.0,
  ↪ 'O': 0.0}
4 And, the optimal policy is:
5 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
  ↪ 'O': 'r'}
```

Problem 3: $Q^*(A, \text{down})$ and $Q^*(A, \text{right})$ with $\gamma=1$

Assume all the non-exit actions with a success probability of p that it transitions to next state and acquires the reward along the edge, otherwise, the agent would state in the same state with a reward of 0. If the discount factor $\gamma=1$, please calculate the $Q^*(A, \text{down})$ and $Q^*(A, \text{right})$, where action 'down' refers that the agent move from A to H, while action 'right' refers that the agent moves from A to B. **Answer:**

For answering this question, I have Q^* for all State-Actions as well as A and it's possible actions (e.g., right and down). I see that value iteration with $\gamma=1$ and $p = 0.5$ will be converged in 55 iterations. If we increase the probability to 0.99, it takes 3235 iterations to be converged, and if we set $p = 1$, it won't be converged, which is interesting.

The output of my code for this problem is presented below:

Code output for Problem 3

```

1 Q*( A ) for actions right, left, up, down, and exit are = dict_values([1.0, 0.0,
  ↪ 0.0, 1.0, 0.0])
2 After 55 iterations:
3 The final value for each state is:
4 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
  ↪ 'O': 0.0}
5 And, the optimal policy is:
6 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
  ↪ 'O': 'r'}
```

```

7
8 Q*( B ) for actions right, left, up, down, and exit are = dict_values([-1.125, 1.0,
  ↪ 0.0, 0.0, 0.0])
9 After 55 iterations:
10 The final value for each state is:
11 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
  ↪ 'O': 0.0}
12 And, the optimal policy is:
13 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
  ↪ 'O': 'r'}
```

```

14
15 Q*( C ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
  ↪ 0.0, 7.75, 0.0])
16 After 55 iterations:
17 The final value for each state is:
18 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
  ↪ 'O': 0.0}
19 And, the optimal policy is:
20 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
  ↪ 'O': 'r'}
```

```

21
22 Q*( D ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
  ↪ 0.0, 25.5, 0.0])
23 After 55 iterations:
```

```

24 The final value for each state is:
25 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
26 And, the optimal policy is:
27 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
28
29 Q*( E ) for actions right, left, up, down, and exit are = dict_values([0.0, 1.0,
    ↪ 0.0, 0.0, 0.0])
30 After 55 iterations:
31 The final value for each state is:
32 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
33 And, the optimal policy is:
34 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
35
36 Q*( F ) for actions right, left, up, down, and exit are = dict_values([1.0, 1.0,
    ↪ 0.0, 0.0, 0.0])
37 After 55 iterations:
38 The final value for each state is:
39 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
40 And, the optimal policy is:
41 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
42
43 Q*( G ) for actions right, left, up, down, and exit are = dict_values([1.0, 0.0,
    ↪ 0.0, 0.0, 0.0])
44 After 55 iterations:
45 The final value for each state is:
46 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
47 And, the optimal policy is:
48 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
49
50 Q*( H ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
    ↪ 1.0, 1.0, 0.0])
51 After 55 iterations:
52 The final value for each state is:
53 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
54 And, the optimal policy is:
55 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
56
57 Q*( O ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
    ↪ 0.0, 0.0, 0.0])
58 After 55 iterations:
59 The final value for each state is:
60 {'A': 1.0, 'B': 1.0, 'C': 7.75, 'D': 25.5, 'E': 1.0, 'F': 1.0, 'G': 1.0, 'H': 1.0,
    ↪ 'O': 0.0}
61 And, the optimal policy is:
62 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}

```

Problem 4: $Q^*(A, \text{down})$ and $Q^*(A, \text{right})$ with $\gamma=0.5$

Assume all the non-exit actions with a success probability of p that it transitions to next state and acquires the reward along the edge, otherwise, the agent would state in the same state with a reward of 0. If the discount factor $\gamma=0.5$, please calculate the $Q^*(A, \text{down})$ and $Q^*(A, \text{right})$, where action ‘down’ refers that the agent move from A to H, while action ‘right’ refers that the agent moves from A to B.

Answer:

And the output of the code for value iteration with $\gamma=0.5$ and $p = 0.5$ will be:

Code output for Problem 3

```

1  Q*( A ) for actions right, left, up, down, and exit are =
   ↪ dict_values([0.6666666666666666, 0.0, 0.0, 0.6666666666666666, 0.0])
2  After 28 iterations:
3  The final value for each state is:
4  {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
   ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
   ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
5  And, the optimal policy is:
6  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
   ↪ 'O': 'r'}
7
8  Q*( B ) for actions right, left, up, down, and exit are =
   ↪ dict_values([-4.677083333333333, 0.6666666666666666, 0.0, 0.0, 0.0])
9  After 28 iterations:
10 The final value for each state is:
11 {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
   ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
   ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
12 And, the optimal policy is:
13 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
   ↪ 'O': 'r'}
14
15 Q*( C ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
   ↪ 0.0, 1.2916666666666667, 0.0])
16 After 28 iterations:
17 The final value for each state is:
18 {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
   ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
   ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
19 And, the optimal policy is:
20 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
   ↪ 'O': 'r'}
21
22 Q*( D ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
   ↪ 0.0, 25.166666666666668, 0.0])
23 After 28 iterations:
24 The final value for each state is:
25 {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
   ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
   ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
26 And, the optimal policy is:
27 {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
   ↪ 'O': 'r'}
28
29 Q*( E ) for actions right, left, up, down, and exit are = dict_values([0.0,
   ↪ 0.6666666666666666, 0.0, 0.0, 0.0])
30 After 28 iterations:
31 The final value for each state is:
32 {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
   ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
   ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
33 And, the optimal policy is:

```

```

34  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
35
36  Q*( F ) for actions right, left, up, down, and exit are =
    ↪ dict_values([0.6666666666666666, 0.6666666666666666, 0.0, 0.0, 0.0])
37  After 28 iterations:
38  The final value for each state is:
39  {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
    ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
    ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
40  And, the optimal policy is:
41  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
42
43  Q*( G ) for actions right, left, up, down, and exit are =
    ↪ dict_values([0.6666666666666666, 0.0, 0.0, 0.0, 0.0])
44  After 28 iterations:
45  The final value for each state is:
46  {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
    ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
    ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
47  And, the optimal policy is:
48  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
49
50  Q*( H ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
    ↪ 0.6666666666666666, 0.6666666666666666, 0.0])
51  After 28 iterations:
52  The final value for each state is:
53  {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
    ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
    ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
54  And, the optimal policy is:
55  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}
56
57  Q*( O ) for actions right, left, up, down, and exit are = dict_values([0.0, 0.0,
    ↪ 0.0, 0.0, 0.0])
58  After 28 iterations:
59  The final value for each state is:
60  {'A': 0.6666666666666666, 'B': 0.6666666666666666, 'C': 1.2916666666666667, 'D':
    ↪ 25.166666666666668, 'E': 0.6666666666666666, 'F': 0.6666666666666666, 'G':
    ↪ 0.6666666666666666, 'H': 0.6666666666666666, 'O': 0.0}
61  And, the optimal policy is:
62  {'A': 'r', 'B': 'l', 'C': 'd', 'D': 'd', 'E': 'l', 'F': 'r', 'G': 'r', 'H': 'u',
    ↪ 'O': 'r'}

```

Appendix I

Value iteration code for the desired MDP problem.

Value Iteration Code

```

1  #!/usr/bin/env python
2
3  '''Pourya Shahverdi
4      CSI-5130 Artificial Intelligence
5      MDP
6      '''
7
8  import sys
9
10 p = 1 #Probabiliy
11 GAMMA = 0.5 #Discount Factor
12 x = 1 #Reward on the edge from H to G (i.e., H ---- x ----> G)
13
14 all_states = ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'O']
15 all_actions = ['r', 'l', 'u', 'd', 'e']
16
17 #Defining the Environment: environment(s, a) receives state-action pairs as
18 # input and returns their probabilities and rewards for the corresponding successor.
19 def environment(s, a):
20     if s == 'A' and a == 'd':
21         return p, 1, 'H'
22     elif s == 'A' and a == 'd':
23         return 1-p, 0, 'A'
24     elif s == 'A' and a == 'r':
25         return p, 1, 'B'
26     elif s == 'A' and a == 'r':
27         return 1-p, 0, 'A'
28     elif s == 'B' and a == 'l':
29         return p, 1, 'A'
30     elif s == 'B' and a == 'l':
31         return 1-p, 0, 'B'
32     elif s == 'B' and a == 'r':
33         return p, -10, 'C'
34     elif s == 'B' and a == 'r':
35         return 1-p, 0, 'B'
36     elif s == 'C' and a == 'd':
37         return p, -10, 'D'
38     elif s == 'C' and a == 'd':
39         return 1-p, 0, 'C'
40     elif s == 'D' and a == 'd':
41         return p, 50, 'E'
42     elif s == 'D' and a == 'd':
43         return 1-p, 0, 'D'
44     elif s == 'E' and a == 'l':
45         return p, 1, 'F'
46     elif s == 'E' and a == 'l':
47         return 1-p, 0, 'E'
48     elif s == 'F' and a == 'r':
49         return p, 1, 'E'
50     elif s == 'F' and a == 'r':
51         return 1-p, 0, 'F'
52     elif s == 'F' and a == 'l':
53         return p, 1, 'G'
54     elif s == 'F' and a == 'l':
55         return 1-p, 0, 'F'
56     elif s == 'G' and a == 'r':
57         return p, 1, 'F'
58     elif s == 'G' and a == 'r':

```

```

59     return 1-p, 0, 'G'
60 elif s == 'H' and a == 'd':
61     return p, x, 'G'
62 elif s == 'H' and a == 'd':
63     return 1-p, 0, 'H'
64 elif s == 'H' and a == 'u':
65     return p, 1, 'A'
66 elif s == 'H' and a == 'u':
67     return 1-p, 0, 'H'
68 else:
69     return p, 0, 'Q'
70
71 # Value iteration function:
72 def value_iteration(all_states, all_actions, s_a_p_r_suc):
73     V = {s: 0 for s in all_states}
74     k = 0
75     optimal_policy = {s: 0 for s in all_states}
76     f = open('value_iteration.txt', "w") #To clear the txt file.
77     f.close()
78     while True:
79         k = k+1
80         oldV = V.copy()
81         for s in all_states:
82             Q = {}
83             for a in all_actions:
84                 for suc in environment(s,a)[2]:
85                     Q[a] = environment(s,a)[0]*(environment(s,a)[1] + GAMMA *
86                         ↪ oldV[suc])
87             # print("Q*(",s,") for actions right, left, up, down, and exit are
88             ↪ "=",Q.values())
89             V[s] = max(Q.values())
90             optimal_policy[s] = max(Q, key=Q.get)
91             # print("Q*[A] = ", V['A'],"\n")
92             # print("Values of all states at the end of iteration #",k,"= \n", V)
93             # print("\n")
94             f = open('value_iteration.txt', 'a')
95             f.seek(0)
96             data = f.read()
97             if len(data) > 0:
98                 f.write("\n")
99             print(V,optimal_policy, file = f)
100             f.close()
101             if all(oldV[s] == V[s] for s in all_states):
102                 break
103         return V,k,optimal_policy
104
105 if __name__ == "__main__":
106     print("After",value_iteration(all_states, all_actions,
107         ↪ environment)[1],"iterations:")
108     print("The final value for each state is:\n",value_iteration(all_states,
109         ↪ all_actions, environment)[0])
110     print("And, the optimal policy is:\n", value_iteration(all_states, all_actions,
111         ↪ environment)[2])
112
113     f = open('value_iteration.txt', 'a')
114     f.seek(0)
115     data = f.read()
116     if len(data) > 0:
117         f.write("\n")
118     print("After",value_iteration(all_states, all_actions,
119         ↪ environment)[1],"iterations:",file=f)
120     print("The final value for each state is:\n",value_iteration(all_states,
121         ↪ all_actions, environment)[0],file=f)

```

```
116 print("And, the optimal policy is:\n", value_iteration(all_states, all_actions,  
    ↪ environment)[2],file=f)
```