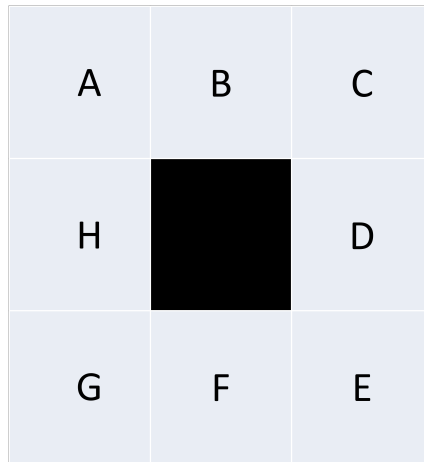


## § Reinforcement Learning §



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

Given 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 (right, left, up, down) and exit. An agent acting in this environment has recorded the following episodes:

### Episode 1

s	a	s'	r
A	right	B	1
B	Right	C	-10
C	Down	D	-10
D	Down	E	50

### Episode 2

s	A	s'	r
A	Right	B	1
B	Left	A	1
A	Right	B	1
B	Right	C	-10
C	Down	D	-10
D	Down	E	50

## Episode 3

s	A	s'	r
A	down	H	1
H	down	G	1
G	Right	F	1
F	Right	E	1

## Episode 4

s	A	s'	R
A	down	H	1
H	up	A	1
A	down	H	1
H	down	G	1
G	up	H	1
H	down	G	1
G	Right	F	1
F	Left	G	1
G	Right	F	1
F	Right	E	1

Codes are publicly available at [GitHub: RL-Problem-AI-CSI-5130](https://github.com/PouryaShahverdi/RL-Problem-AI-CSI-5130)

## Problem 1: Model-Based RL, Direct Evaluation

Consider the agent running model-based reinforcement learning (direct evaluation) based on the four episodes collected above, calculate the transition probabilities, and reward functions.

**Answer:**

*Model-Based*  
1- Direct Evaluation:

s	a	s'	count	$T(s,a,s')$	$R(s,a,s')$
A	r	B	3	$\frac{3}{5} \times 1$	1
A	d	H	3	$\frac{3}{5} \times 1$	1
B	r	C	2	$\frac{2}{2} \times 1$	-10
B	d	A	1	$\frac{1}{1} \times 1$	1
C	d	D	2	$\frac{2}{2} \times 1$	-10
D	d	E	2	$\frac{2}{2} \times 1$	50
F	r	F	2	$\frac{2}{2} \times 1$	1
F	d	G	1	$\frac{1}{1} \times 1$	1
G	r	F	3	$\frac{3}{3} \times 1$	1
G	u	H	1	$\frac{1}{1} \times 1$	1
H	d	G	3	$\frac{3}{3} \times 1$	1
H	u	A	1	$\frac{1}{1} \times 1$	1

## Problem 2: Model-Based RL, Temporal Difference

Consider the agent running model-based reinforcement learning (temporal difference) based on the four episodes collected above, calculate the values of each state. We assume the discount factor  $\gamma = 1/2$ , and learning rate  $\alpha = 1/2$ .

**Answer:**

I've written a code to solve this problem as follows:

Value Iteration Code

```

1  #!/usr/bin/env python
2  from math import gamma
3
4  alpha = 0.5
5  gamma = 0.5
6
7  all_states = ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'O']
8  all_actions = ['r', 'l', 'u', 'd', 'e']
9
10 # episode1 = [['A', 'r', 'B', 1],
11 #              ['B', 'r', 'C', -10],
12 #              ['C', 'd', 'D', -10],
13 #              ['D', 'd', 'E', 50]]
14
15 # episode2 = [['A', 'r', 'B', 1],
16 #              ['B', 'l', 'A', 1],
17 #              ['A', 'r', 'B', 1],
18 #              ['B', 'r', 'C', -10],
19 #              ['C', 'd', 'D', -10],
20 #              ['D', 'd', 'E', 50]]
21
22 # episode3 = [['A', 'd', 'H', 1],
23 #              ['H', 'd', 'G', 1],
24 #              ['G', 'r', 'F', 1],
25 #              ['F', 'r', 'E', 1]]
26
27 # episode4 = [['A', 'd', 'H', 1],
28 #              ['H', 'u', 'A', 1],
29 #              ['A', 'd', 'H', 1],
30 #              ['H', 'd', 'G', 1],
31 #              ['G', 'u', 'H', 1],
32 #              ['H', 'd', 'G', 1],
33 #              ['G', 'r', 'F', 1],
34 #              ['F', 'l', 'G', 1],
35 #              ['G', 'r', 'F', 1],
36 #              ['F', 'r', 'E', 1]]
37
38 episodes = [[['A', 'r', 'B', 1],
39               ['B', 'r', 'C', -10],
40               ['C', 'd', 'D', -10],
41               ['D', 'd', 'E', 50]],
42              [['A', 'r', 'B', 1],
43               ['B', 'l', 'A', 1],
44               ['A', 'r', 'B', 1],
45               ['B', 'r', 'C', -10],
46               ['C', 'd', 'D', -10],
47               ['D', 'd', 'E', 50]],
48              [['A', 'd', 'H', 1],
49               ['H', 'd', 'G', 1],
50               ['G', 'r', 'F', 1],
51               ['F', 'r', 'E', 1]],
52              [['A', 'd', 'H', 1],

```

```

53         ['H', 'u', 'A', 1],
54         ['A', 'd', 'H', 1],
55         ['H', 'd', 'G', 1],
56         ['G', 'u', 'H', 1],
57         ['H', 'd', 'G', 1],
58         ['G', 'r', 'F', 1],
59         ['F', 'l', 'G', 1],
60         ['G', 'r', 'F', 1],
61         ['F', 'r', 'E', 1]]]
62
63 V = {state : 0 for state in all_states}
64
65 k = 0
66 for episode in episodes:
67     for samples in episode:
68         oldV = V.copy()
69         for state in samples[0]:
70             for action in samples[1]:
71                 for successor in samples[2]:
72                     reward = samples[3]
73                     V[state] = (1-alpha)*oldV[state] + (alpha * (reward + gamma *
74                                     ↳ V[successor]))
75                     print('Value of States:',V)
76                     print('Value of',successor,'as the successor
77                             ↳ for',state,'is:',V[successor])

```

And the result is:

Code output for Problem 2

```

1 ('Value of States:', {'A': 0.5, 'C': 0, 'B': 0, 'E': 0, 'D': 0, 'G': 0, 'F': 0, 'H':
  ↳ 0, 'O': 0})
2 ('Value of', 'B', 'as the successor for', 'A', 'is:', 0)
3 ('Value of States:', {'A': 0.5, 'C': 0, 'B': -5.0, 'E': 0, 'D': 0, 'G': 0, 'F': 0,
  ↳ 'H': 0, 'O': 0})
4 ('Value of', 'C', 'as the successor for', 'B', 'is:', 0)
5 ('Value of States:', {'A': 0.5, 'C': -5.0, 'B': -5.0, 'E': 0, 'D': 0, 'G': 0, 'F':
  ↳ 0, 'H': 0, 'O': 0})
6 ('Value of', 'D', 'as the successor for', 'C', 'is:', 0)
7 ('Value of States:', {'A': 0.5, 'C': -5.0, 'B': -5.0, 'E': 0, 'D': 25.0, 'G': 0,
  ↳ 'F': 0, 'H': 0, 'O': 0})
8 ('Value of', 'E', 'as the successor for', 'D', 'is:', 0)
9 ('Value of States:', {'A': -0.5, 'C': -5.0, 'B': -5.0, 'E': 0, 'D': 25.0, 'G': 0,
  ↳ 'F': 0, 'H': 0, 'O': 0})
10 ('Value of', 'B', 'as the successor for', 'A', 'is:', -5.0)
11 ('Value of States:', {'A': -0.5, 'C': -5.0, 'B': -2.125, 'E': 0, 'D': 25.0, 'G': 0,
  ↳ 'F': 0, 'H': 0, 'O': 0})
12 ('Value of', 'A', 'as the successor for', 'B', 'is:', -0.5)
13 ('Value of States:', {'A': -0.28125, 'C': -5.0, 'B': -2.125, 'E': 0, 'D': 25.0, 'G':
  ↳ 0, 'F': 0, 'H': 0, 'O': 0})
14 ('Value of', 'B', 'as the successor for', 'A', 'is:', -2.125)
15 ('Value of States:', {'A': -0.28125, 'C': -5.0, 'B': -7.3125, 'E': 0, 'D': 25.0,
  ↳ 'G': 0, 'F': 0, 'H': 0, 'O': 0})
16 ('Value of', 'C', 'as the successor for', 'B', 'is:', -5.0)
17 ('Value of States:', {'A': -0.28125, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 25.0,
  ↳ 'G': 0, 'F': 0, 'H': 0, 'O': 0})
18 ('Value of', 'D', 'as the successor for', 'C', 'is:', 25.0)
19 ('Value of States:', {'A': -0.28125, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
  ↳ 'G': 0, 'F': 0, 'H': 0, 'O': 0})
20 ('Value of', 'E', 'as the successor for', 'D', 'is:', 0)
21 ('Value of States:', {'A': 0.359375, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
  ↳ 'G': 0, 'F': 0, 'H': 0, 'O': 0})
22 ('Value of', 'H', 'as the successor for', 'A', 'is:', 0)
23 ('Value of States:', {'A': 0.359375, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,

```

```

    → 'G': 0, 'F': 0, 'H': 0.5, 'O': 0})
24 ('Value of', 'G', 'as the successor for', 'H', 'is:', 0)
25 ('Value of States:', {'A': 0.359375, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
    → 'G': 0.5, 'F': 0, 'H': 0.5, 'O': 0})
26 ('Value of', 'F', 'as the successor for', 'G', 'is:', 0)
27 ('Value of States:', {'A': 0.359375, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
    → 'G': 0.5, 'F': 0.5, 'H': 0.5, 'O': 0})
28 ('Value of', 'E', 'as the successor for', 'F', 'is:', 0)
29 ('Value of States:', {'A': 0.8046875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
    → 'G': 0.5, 'F': 0.5, 'H': 0.5, 'O': 0})
30 ('Value of', 'H', 'as the successor for', 'A', 'is:', 0.5)
31 ('Value of States:', {'A': 0.8046875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D': 37.5,
    → 'G': 0.5, 'F': 0.5, 'H': 0.951171875, 'O': 0})
32 ('Value of', 'A', 'as the successor for', 'H', 'is:', 0.8046875)
33 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 0.5, 'F': 0.5, 'H': 0.951171875, 'O': 0})
34 ('Value of', 'H', 'as the successor for', 'A', 'is:', 0.951171875)
35 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 0.5, 'F': 0.5, 'H': 1.1005859375, 'O': 0})
36 ('Value of', 'G', 'as the successor for', 'H', 'is:', 0.5)
37 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.025146484375, 'F': 0.5, 'H': 1.1005859375, 'O': 0})
38 ('Value of', 'H', 'as the successor for', 'G', 'is:', 1.1005859375)
39 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.025146484375, 'F': 0.5, 'H': 1.30657958984375, 'O': 0})
40 ('Value of', 'G', 'as the successor for', 'H', 'is:', 1.025146484375)
41 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.1375732421875, 'F': 0.5, 'H': 1.30657958984375, 'O': 0})
42 ('Value of', 'F', 'as the successor for', 'G', 'is:', 0.5)
43 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.1375732421875, 'F': 1.034393310546875, 'H': 1.30657958984375, 'O':
    → 0})
44 ('Value of', 'G', 'as the successor for', 'F', 'is:', 1.1375732421875)
45 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.3273849487304688, 'F': 1.034393310546875, 'H': 1.30657958984375,
    → 'O': 0})
46 ('Value of', 'F', 'as the successor for', 'G', 'is:', 1.034393310546875)
47 ('Value of States:', {'A': 1.14013671875, 'C': -1.25, 'B': -7.3125, 'E': 0, 'D':
    → 37.5, 'G': 1.3273849487304688, 'F': 1.0171966552734375, 'H': 1.30657958984375,
    → 'O': 0})
48 ('Value of', 'E', 'as the successor for', 'F', 'is:', 0)

```

## Problem 3: Model-Free RL, Q-Learning

Consider the agent running model-free reinforcement learning (Q-learning) based on the four episodes being repeatedly in an infinite sequence E1, E2, E3, E4, E1, E2, E3, E4...., calculate the Q-values of each <state, action>. We assume the discount factor  $\gamma = 1$ , and learning rate  $\alpha = 1/2$ . [You are encouraged to compute this question with programming assistance.]

### Answer:

I've written a code for this problem as follows:

#### Value Iteration Code

```

1  #!/usr/bin/env python
2  from math import gamma
3
4
5  alpha = 0.5
6  gamma = 1

```

```

7
8 all_states = ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'O']
9 all_actions = ['r', 'l', 'u', 'd', 'e']
10
11 episodes = [[['A', 'r', 'B', 1],
12               ['B', 'r', 'C', -10],
13               ['C', 'd', 'D', -10],
14               ['D', 'd', 'E', 50]],
15              [['A', 'r', 'B', 1],
16               ['B', 'l', 'A', 1],
17               ['A', 'r', 'B', 1],
18               ['B', 'r', 'C', -10],
19               ['C', 'd', 'D', -10],
20               ['D', 'd', 'E', 50]],
21              [['A', 'd', 'H', 1],
22               ['H', 'd', 'G', 1],
23               ['G', 'r', 'F', 1],
24               ['F', 'r', 'E', 1]],
25              [['A', 'd', 'H', 1],
26               ['H', 'u', 'A', 1],
27               ['A', 'd', 'H', 1],
28               ['H', 'd', 'G', 1],
29               ['G', 'u', 'H', 1],
30               ['H', 'd', 'G', 1],
31               ['G', 'r', 'F', 1],
32               ['F', 'l', 'G', 1],
33               ['G', 'r', 'F', 1],
34               ['F', 'r', 'E', 1]]]
35
36 Qmax = {}
37 Q = {}
38 for state in all_states:
39     for action in all_actions:
40         Qmax[state,action] = 0
41         Q[state,action] = 0
42
43 k = 0
44 while k < 100:
45     for episode in episodes:
46         for samples in episode:
47             Old_Q = Q.copy()
48             for state in samples[0]:
49                 for action in samples[1]:
50                     for successor in samples[2]:
51                         reward = samples[3]
52                         Q[state,action] = (1-alpha)*Old_Q[state,action] +
53                             ↪ alpha*(reward + gamma * Qmax[successor,action])
54             for act in all_actions:
55                 if Q[state,act] > Qmax[state,action]:
56                     Qmax[state,action] = Q[state,act]
57
58             print(Q)
59
60     k += 1

```

And the result after 100 iteration is:

Code output for Problem 3

```

1 {('F', 'u'): 0, ('H', 'l'): 0, ('C', 'r'): 0, ('D', 'd'): 50.0, ('E', 'r'): 0, ('E',
  ↪ 'e'): 0, ('B', 'u'): 0, ('E', 'l'): 0, ('F', 'r'): 1.0, ('C', 'l'): 0, ('A',
  ↪ 'e'): 0, ('H', 'd'): 1.0, ('H', 'r'): 0, ('G', 'e'): 0, ('D', 'e'): 0, ('O',
  ↪ 'd'): 0, ('B', 'l'): 1.0, ('E', 'd'): 0, ('C', 'u'): 0, ('F', 'd'): 0, ('A',
  ↪ 'd'): 2.0, ('B', 'r'): -10.0, ('H', 'e'): 0, ('O', 'u'): 0, ('E', 'u'): 0, ('G',
  ↪ 'd'): 0, ('B', 'd'): 0, ('G', 'r'): 2.0, ('A', 'r'): 2.0, ('F', 'l'): 1.0, ('A',
  ↪ 'l'): 0, ('F', 'e'): 0, ('G', 'u'): 2.0, ('D', 'u'): 0, ('G', 'l'): 0, ('B',

```

```
→ 'e'): 0, ('C', 'e'): 0, ('H', 'u'): 1.0, ('A', 'u'): 0, ('O', 'e'): 0, ('D',  
→ 'l'): 0, ('O', 'l'): 0, ('D', 'r'): 0, ('O', 'r'): 0, ('C', 'd'): 40.0}
```

Looks like numbers will be converged after like 54 iterations.