

Assignment 2

Peng Tang 517020910038

In this assignment, we implement first-visit and every-visit MC method and TD(0) to evaluate the random policy on the given small gridworld.

1 Analysis

1.1 First-visit MC

In order to reduce memory use, we use increment Monte-Carlo to update the value of every states. We implement first-visit MC as following steps.

- Producing a episode following current policy π .

Because current policy π is random policy, we randomly choose a state as the start state, then randomly take actions at each state until reaching the terminal states. We record all states we reach in the process.

- After gaining a episode, we traverse the episode and update the value of corresponding states by using the following equations. And for first-visit MC, we just update a state when we visit it for the first time in the episode.

$$G(s_t) = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1-t} R_n = \begin{cases} (-1) \times (n - t) & \gamma = 1 \\ (-1) \times \frac{(1 - \gamma^{n-t})}{1 - \gamma} & 0 < \gamma < 1 \end{cases}$$
$$N(s_t) = N(s_t) + 1$$
$$V(s_t) = V(s_t) + \frac{1}{N(s_t)} (G(s_t) - V(s_t))$$

Where s_t means the t -th state in the episode and n means the number of states in the episode. And for this problem, we have $R_{t+1} = \dots = R_n = -1$. γ is the discount factor.

We do above two steps for enough times, then the values of all states will converge to fixed values.

1.2 Every-visit MC

The difference between every-visit MC and first-visit MC is that every-visit MC need update a state every time we visit it. So we can implement this just like first-visit MC.

1.3 TD(0)

For TD(0), we update the value of states for every step rather than producing a episode. We can do this as following steps.

- We randomly choose a start state and take random actions until reaching a terminal state.
- In the process, everytime we take an action, we can update the value of current state by using following equations.

$$V(s) = V(s) + \alpha(G - V(s)) = V(s) + \alpha(r + \gamma V(s') - V(s))$$

Where s' is the next state of s after taking an action and $r = -1$.

- We do above two steps for enough times and the values of all states will converge to fixed values.

2 Code

In the *solution.py*, we implement MC method and TD(0) method by defining two class, **MC** and **TD**.

2.1 MC Class

```
1  class MC(object):
2
3      def __init__(self, discount):
4
5          self.FirstValue = np.random.randn(6,6) #初始化所有位置的值
6          self.EveryValue = np.random.randn(6,6) #初始化所有位置的值
7          self.FirstValue[0][1] = 0
8          self.FirstValue[5][5] = 0
9          self.EveryValue[0][1] = 0
10         self.EveryValue[5][5] = 0
11
12
13         self.reward = -1
14         self.discount = discount
15         self.action = [[-1,0],[0,1],[1,0],[0,-1]]
16         self.terminalState = np.array([[0,1],[5,5]])
17
18
19     def produce_random_episodes(self):
20         episodesList = []
21
22         state = np.random.randint(0,6,2)
23         episodesList.append(state.copy())
24         while (state!=self.terminalState[0]).any() and
25         (state!=self.terminalState[1]).any():
26             a = np.random.randint(0,4,1)
27             state[0] = max(0,min(5,state[0]+self.action[a[0]][0]))
28             state[1] = max(0,min(5,state[1]+self.action[a[0]][1]))
29
30             episodesList.append(state.copy())
31
32         return episodesList
33
34     def first_visit_MC(self,maxEpisodes):
35         count = np.zeros((6,6))
36
37         times = 1
38         error = []
39         while times < maxEpisodes:
40
41             tempcount = np.zeros((6,6))
42             episod = self.produce_random_episodes()
43
44             n = len(episod)
```

```

44
45         maxerror = 0
46         for i in range(n-1):
47
48             if tempcount[episd[i][0]][episd[i][1]] == 0:
49                 tempcount[episd[i][0]][episd[i][1]] = 1
50                 count[episd[i][0]][episd[i][1]] +=1
51                 Gt = 0
52                 if self.discount == 1:
53                     Gt = self.reward*(n-i-1)
54                 else:
55                     Gt = self.reward*((1-self.discount**(n-i-1))/(1-
self.discount))
56
57                 oldval = self.FirstValue[episd[i][0]][episd[i][1]]
58                 self.FirstValue[episd[i][0]][episd[i][1]] +=
1./count[episd[i][0]][episd[i][1]]*(Gt-self.FirstValue[episd[i][0]][episd[i]
[1]])
59
60                 maxerror = max(maxerror,abs(oldval-
self.FirstValue[episd[i][0]][episd[i][1]]))
61                 times+=1
62                 error.append(maxerror)
63
64         return error
65
66     def every_visit_MC(self,maxEpisodes):
67
68         count = np.zeros((6,6))
69
70         times = 1
71         error = []
72         while times < maxEpisodes:
73
74             episod = self.produce_random_episodes()
75             n = len(episod)
76             maxerror = 0
77             for i in range(n-1):
78                 count[episd[i][0]][episd[i][1]] +=1
79                 Gt = 0
80                 if self.discount == 1:
81                     Gt = self.reward*(n-i-1)
82                 else:
83                     Gt = self.reward*((1-self.discount**(n-i-1))/(1-
self.discount))
84
85                 oldval = self.EveryValue[episd[i][0]][episd[i][1]]
86                 self.EveryValue[episd[i][0]][episd[i][1]] +=
1./count[episd[i][0]][episd[i][1]]*(Gt-self.EveryValue[episd[i][0]][episd[i]
[1]])
87
88                 maxerror = max(maxerror,abs(oldval-
self.EveryValue[episd[i][0]][episd[i][1]]))
89                 times+=1
90                 error.append(maxerror)
91
92         return error

```

- In the class, **FirstValue** is used to store the values for first-visit method; **EveryValue** is used to store the values for every-visit method. **reward** is always -1 in this problem; **discount** is

the discount factor; **action** contains four actions that we can take; **terminalState** contains two terminal states.

- The **produce_random_episodes()** function is used to produce a episode randomly and return a list containing all states of the episode.
- The **first_visit_MC()** function is used to compute the values with first-visit MC method. In this function, we need a parameter as the terminal condition. We implement the function as the analysis.
- The **every_visit_MC()** function is used to compute the values with every-visit MC method. In this function, we need a parameter as the terminal condition. We implement the function as the analysis.

2.2 TD Class

```
1 class TD(object):
2     def __init__(self, discount, stepSize):
3
4         self.Value = np.random.randn(6,6)
5
6         self.Value[0][1] = 0
7         self.Value[5][5] = 0
8         self.reward = -1
9         self.discount = discount
10        self.stepSize = stepSize
11        self.action = [[-1,0],[0,1],[1,0],[0,-1]]
12        self.terminalState = np.array([[0,1],[5,5]])
13
14
15    def TD_find(self, maxepisodes):
16
17        times = 1
18        error = []
19        while times < maxepisodes:
20            self.stepSize = self.stepSize*0.99
21            maxerror = 0
22            state = np.random.randint(0,6,2)
23            while (state!=self.terminalState[0]).any() and
24            (state!=self.terminalState[1]).any():
25                a = np.random.randint(0,4,1)
26
27                oldval = self.Value[state[0]][state[1]]
28                oldstate = state.copy()
29                state[0] = max(0,min(5,state[0]+self.action[a[0]][0]))
30                state[1] = max(0,min(5,state[1]+self.action[a[0]][1]))
31                self.Value[oldstate[0]][oldstate[1]] = self.stepSize*
32                (self.reward + self.discount*self.Value[state[0]][state[1]]) + (1-
33                self.stepSize)*self.Value[oldstate[0]][oldstate[1]]
34
35                maxerror = max(maxerror, abs(oldval-self.Value[oldstate[0]]
36                [oldstate[1]]))
37                times+=1
38
39            error.append(maxerror)
40        return error
```

- In this class, the **stepSize** is the weight factor α we use to update the value.

- The **TD_find()** function is used to compute the values with TD(0) method. In this function, we need a parameter as the terminal condition. We implement the function as the analysis.

2.3 Main function

```
1  def main():
2      MC_discount = 1
3      MC_maxEpsides = 5000
4      MC_model = MC(MC_discount)
5      t1 = time.time()
6      Error1 = MC_model.first_visit_MC(MC_maxEpsides)
7      t2 = time.time()
8      Error2 = MC_model.every_visit_MC(MC_maxEpsides)
9      t3 = time.time()
10
11     print("The spent time of MC first visit: %s" %(t2-t1))
12     print("The spent time of MC every visit: %s" %(t3-t2))
13
14
15
16     #TD
17     TD_discount = 1
18     TD_maxEpsides = 5000
19     TD_stepsize = 0.8
20     TD_model = TD(TD_discount,TD_stepsize)
21     t4 = time.time()
22     Error3 = TD_model.TD_find(TD_maxEpsides)
23     t5 = time.time()
24     print("The spent time of TD: %s" %(t5-t4))
25
26
27
28     plotResult(MC_model.FirstValue,MC_model.EveryValue,TD_model.value)
29
30     pltError(Error1,Error2,Error3)
```

This function is used to call the methods we implement and produce the results.

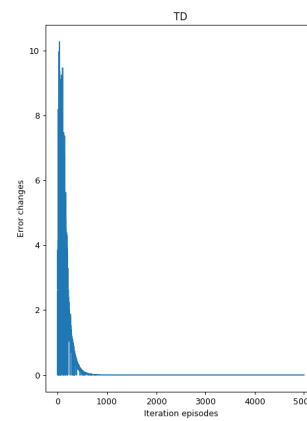
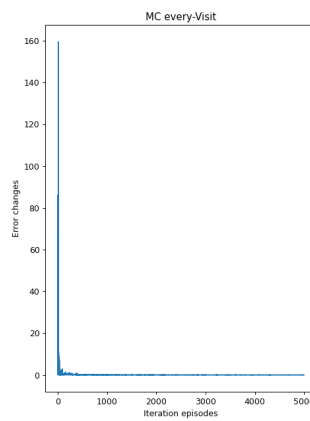
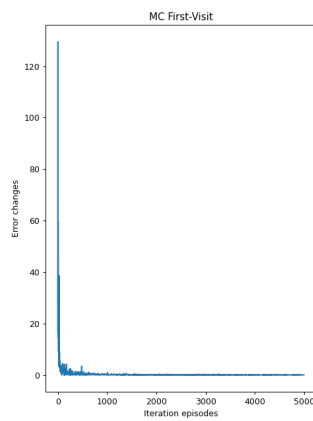
3 Results

By running the code, we can get following results.

| MC First-Visit | | | | | |
|----------------|--------|--------|--------|--------|--------|
| -17.6 | 0.0 | -28.91 | -44.46 | -51.73 | -55.23 |
| -30.78 | -28.92 | -39.45 | -47.61 | -51.83 | -53.15 |
| -44.45 | -44.02 | -48.06 | -50.77 | -50.92 | -50.51 |
| -53.58 | -52.38 | -52.83 | -51.81 | -46.75 | -44.31 |
| -58.91 | -57.74 | -54.88 | -49.07 | -41.09 | -30.8 |
| -61.05 | -59.51 | -55.32 | -48.06 | -31.6 | 0.0 |

| MC every-Visit | | | | | |
|----------------|--------|--------|--------|--------|--------|
| -17.12 | 0.0 | -29.58 | -45.39 | -51.66 | -55.52 |
| -32.3 | -30.11 | -39.03 | -46.55 | -51.73 | -54.45 |
| -45.09 | -44.52 | -47.37 | -49.77 | -50.11 | -50.57 |
| -51.4 | -51.12 | -52.08 | -49.37 | -44.97 | -42.78 |
| -57.51 | -55.29 | -53.31 | -47.92 | -38.76 | -28.84 |
| -61.4 | -58.52 | -53.12 | -45.92 | -30.04 | 0.0 |

| TD | | | | | |
|--------|--------|--------|--------|--------|--------|
| -14.95 | 0.0 | -26.47 | -39.82 | -46.6 | -49.33 |
| -27.25 | -24.03 | -34.53 | -42.92 | -46.93 | -48.23 |
| -39.01 | -39.1 | -42.35 | -45.23 | -45.64 | -44.71 |
| -47.27 | -47.21 | -46.76 | -45.2 | -41.68 | -36.98 |
| -51.83 | -50.79 | -48.25 | -42.73 | -34.39 | -21.74 |
| -53.98 | -52.4 | -48.55 | -40.0 | -25.53 | 0.0 |



```

PS E:\大三下\强化学习\作业\homework2\code> D:\python\python.exe "e:\大三下\强化学习\作业\homework2\code\solution.py"
The spent time of MC first visit: 4.1528167724609375
The spent time of MC every visit: 4.19720721244812
The spent time of TD: 4.327538967132568

```

- From the first figure, we can see that first-visit MC and every-visit MC give similar results, and the result of TD(0) is lightly different from their.
- From the second figure, we can see that the convergence rate of every-visit MC and TD(0) is faster than first-visit MC. And the change of TD's error is smaller.
- From the second figure, we can see that the running time of three methods is almost same.