

Assignment 1

Peng Tang 517020910038

In this assignment, we implement iterative policy evaluation, policy iteration and value iteration on the given small gridworld and compare the performance of policy iteration with value iteration.

1 Analysis

1.1 Policy Iteration

For policy iteration, we use policy evaluation and policy improvement alternatively to find the optimal policy for the given gridworld.

- Policy evaluation: For current policy π , we use the following equation to update the values of all states(except the two terminal states) until the biggest change $\Delta < \theta$, θ is a small number.

$$V_{\pi}(s) = \sum_{s'} p(s'|s, \pi(s))[-1 + \gamma V_{\pi}(s')]$$

Then, $V_{\pi}(s)$ will converge to a fixed value following π for all states.

- Policy improvement: After we get the $V_{\pi}(s)$ under a policy π , we can use the following equation to improve current policy.

$$\pi'(s) = \arg \max_a \sum_{s'} p(s'|s, a)[-1 + \gamma V_{\pi}(s')] \quad a \in \{n, e, s, w\}$$

Then, if there is no any change in actions of all states in this computation, we get the optimal policy.

1.2 Value Iteration

For value iteration, we just update the values for all states by using the following equation until they converge to the final values without considering policy.

$$V(s)_{k+1} = \max_{a \in A} (-1 + \gamma \sum_{s'} p_{ss'}^a V_k(s')) \quad A = \{n, e, s, w\}$$

After getting the final values, we can using the following equation to find the optimal policy.

$$\pi(s) = \arg \max_a \sum_{s'} p(s'|s, a)[-1 + \gamma V(s')] \quad a \in \{n, e, s, w\}$$

2 Code

We implement the policy iteration and value iteration in *solution.py* by defining two class, **PolicyIteration** and **ValueIteration**.

2.1 PolicyIteration Class

```

1 class PolicyIteration(object):
2     def __init__(self, discount):
3
4         self.BasicAction = [[0,1],[1,0],[0,-1],[-1,0]] #actions
5         self.action_convert = ["right", "down", "left", "up"]
6         self.Value = np.zeros([6,6]) #begining values are 0s
7         self.error = 0.01
8         self.discntFactor = discount
9         self.reward = -1
10        self.finalState = [[0,1],[5,5]]
11
12        self.changeserror = []
13        #random policy
14        self.Actions = np.empty([6,6], dtype = object)
15        for i in range(6):
16            for j in range(6):
17
18                if [i,j]==self.finalState[0] or [i,j]==self.finalState[1]:
19                    continue
20
21                self.Actions[i][j] = [0,1,2,3]
22
23
24        def PolicyEvaluation(self): #policy evaluation
25
26            while True:
27                maxerror = 0#the biggest error
28
29                old_values = np.copy(self.Value)
30                #traverse
31                for i in range(6):
32                    for j in range(6):
33
34                        if(i==0 and j==1) or (i==5 and j==5): #ingore the
terminal states
35                            continue
36
37                        action = self.Actions[i][j] #current policy
38                        lenth = len(action)
39                        self.Value[i][j] = 0
40                        for k in range(lenth):
41                            ii = max(0,min(i+self.BasicAction[action[k]][0],5))
#compute the next state
42                            jj = max(0,min(j+self.BasicAction[action[k]][1],5))
43                            self.Value[i][j] += 1./lenth*
(self.reward+self.discntFactor*old_values[ii][jj])
44
45                            maxerror = max(maxerror, abs(old_values[i][j]-
self.Value[i][j]))
46                        self.changeserror.append(maxerror)
47                        if maxerror<self.error:
48                            break
49
50        def PolicyImprove(self): #policy improvement
51            policy_state = True
52
53            for i in range(6):
54                for j in range(6):

```

```

55
56         if [i,j]==self.finalState[0] or [i,j]==self.finalState[1]:
57             continue
58
59         old_action = self.Actions[i][j]
60         action_val = []
61         for a in range(4):
62             ii = max(0,min(i+self.BasicAction[a][0],5))
63             jj = max(0,min(j+self.BasicAction[a][1],5))
64
65         action_val.append(self.reward+self.discntFactor*self.value[ii][jj])
66
67         max_action_val = max(action_val)
68         current_action = [k for k in range(4) if action_val[k] ==
max_action_val]
69         self.Actions[i][j] = current_action
70         if current_action != old_action:
71             policy_state = False
72
73         return policy_state
74
75     #state in [0,35]
76     def findPolicy(self, state):
77
78         if state==1 or state==35: #terminal
79             return ["stop"]
80
81         i = state/6
82         j = state - i*6
83         current_state = [i,j]
84         action = self.Actions[i][j]
85         action_list = []
86
87         while current_state != self.finalState[0] and current_state !=
self.finalState[1]:
88             current_state = current_state + self.BasicAction[action[0]]
89             action_list.append(self.action_convert[action[0]])
90
91         return action_list

```

- This class has three functions and several variables to implement the policy iteration.
- **BasicAction** is the actions space for all states; **action_convert** is used as a mapping of **BasicAction**; **Values** is used to record the values of all states; **error** is used as the termination condition for policy evaluation; **discnFactor** is the discount factor for computing the values; **reward** is -1 in this question; **finalState** is to record the terminal states; **Actions** is to record current actions of all states for current policy and the initial policy is the random policy; **changeserror** is used to record the all Δ .
- **PolicyEvaluation()** is used to implement policy evaluation. In this function, we just traverse all states again and again to update the **Values** by using the **Actions** until the max error is smaller than **error**.
- **PolicyImprove()** is used to implement policy improvement. In this function, we traverse all states to update their actions according to new values in **Values**. If there are no changes in **Actions**, we can stop.
- **findPolicy()** is to find the optimal policy for a given state when we already update the **Actions**.

2.2 ValueIteration Class

```
1 class ValueIteration(object):
2     def __init__(self,discout):
3
4
5         self.BasicAction = [[0,1],[1,0],[0,-1],[-1,0]]
6         self.action_convert = ["right","down","left","up"]
7         self.value = np.zeros([6,6])
8         self.Actions = np.empty([6,6],dtype = object)
9         self.error = 0.01
10        self.discntFactor = discount
11        self.reward = -1
12        self.finalState = [[0,1],[5,5]]
13
14        self.changerrors = []
15
16
17
18    def FindValue(self):
19        while True:
20            maxerror = 0
21            old_values = np.copy(self.value)
22            for i in range(6):
23                for j in range(6):
24
25                    if [i,j]==self.finalState[0] or
[i,j]==self.finalState[1]:
26                        continue
27
28                    action_val_max = -inf
29                    for k in range(4):
30                        ii = max(0,min(i+self.BasicAction[k][0],5))
31                        jj = max(0,min(j+self.BasicAction[k][1],5))
32                        action_val_max = max(action_val_max,
(self.reward+self.discntFactor*old_values[ii][jj]))
33
34                        maxerror = max(maxerror,abs(action_val_max-self.value[i]
[j]))
35
36                        self.value[i][j] = action_val_max
37                    self.changerrors.append(maxerror)
38                    if maxerror<self.error:
39                        break
40
41
42        for i in range(6):
43            for j in range(6):
44                if [i,j]==self.finalState[0] or [i,j]==self.finalState[1]:
45                    continue
46
47                action_val = []
48                for k in range(4):
49                    ii = max(0,min(i+self.BasicAction[k][0],5))
50                    jj = max(0,min(j+self.BasicAction[k][1],5))
```

```

51     action_val.append(self.reward+self.discntFactor*self.value[ii][jj])
52
53     max_action_val = max(action_val)
54     current_action = [k for k in range(4) if action_val[k] ==
max_action_val]
55     self.Actions[i][j] = current_action
56
57
58
59     def findPolicy(self, state):
60
61
62         if state==1 or state==35:
63             return ["stop"]
64         i = state/6
65         j = state - i*6
66         current_state = [i,j]
67         action = self.Actions[i][j]
68         action_list = []
69
70         while current_state != self.finalState[0] and current_state !=
self.finalState[0]:
71             current_state = current_state + self.BasicAction[action[0]]
72             action_list.append(self.action_convert[action[0]])
73
74         return action_list

```

- This class has two functions and several variables to implement the value iteration.
- The variables in this class are almost the same as in the **PolicyIteration**.
- **FindValue()** is used to find the convergence values for all states for the first and then find the optimal policy by using these values.
- **findPolicy()** is to find the optimal policy for a given state when we already update the **Actions**.

2.3 Main function

```

1  def main():
2
3      disfactor = 1
4
5
6      #Policy Iteration
7      PolicyMode = PolicyIteration(disfactor)
8
9      flag = False
10     t1 = time.time()
11     while not flag:
12         PolicyMode.PolicyEvaluation()
13         flag = PolicyMode.PolicyImprove()
14     t2 = time.time()
15     PolicyTime = t2-t1

```

```

16 plotResult(PolicyMode.Actions,PolicyMode.Value,"imgs/Policy_Iteration_result.png")
17
18
19
20 #Value Iteration
21
22 ValueMode = ValueIteration(disfactor)
23 t1 = time.time()
24 ValueMode.FindValue()
25 t2 = time.time()
26 ValueTime = t2-t1
27
28
29 plotResult(ValueMode.Actions,ValueMode.Value,"imgs/Value_Iteration_result.png")
30
31 #compare performance
32 plotPerformanceCompare(PolicyMode.changeserror,ValueMode.changerrors)
33 print("The spent time of Policy Iteration: %s" %PolicyTime)
34 print("The spent time of Value Iteration: %s" %ValueTime)

```

By using two class, we can solve this problem with two methods.

3 Results

By running the code, we can get following results.

- Policy Iteration

-1.0	0.0	-1.0	-2.0	-3.0	-4.0
-2.0	-1.0	-2.0	-3.0	-4.0	-4.0
-3.0	-2.0	-3.0	-4.0	-4.0	-3.0
-4.0	-3.0	-4.0	-4.0	-3.0	-2.0
-5.0	-4.0	-4.0	-3.0	-2.0	-1.0
-5.0	-4.0	-3.0	-2.0	-1.0	0.0

→	terminal	←	←	←	←
↖	↑	↖	↖	↖	↓
↖	↑	↖	↖	↘	↓
↖	↑	↖	↘	↘	↓
↖	↑	↘	↘	↘	↓
→	→	→	→	→	terminal

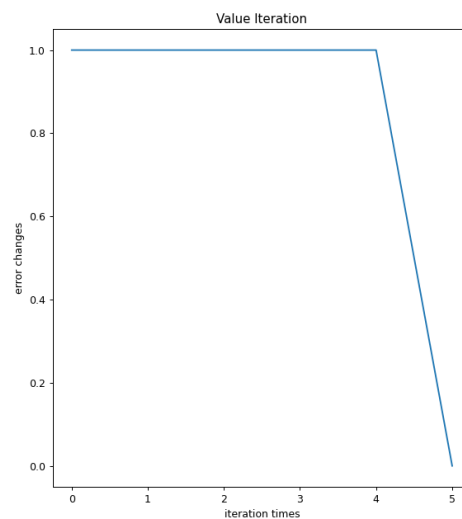
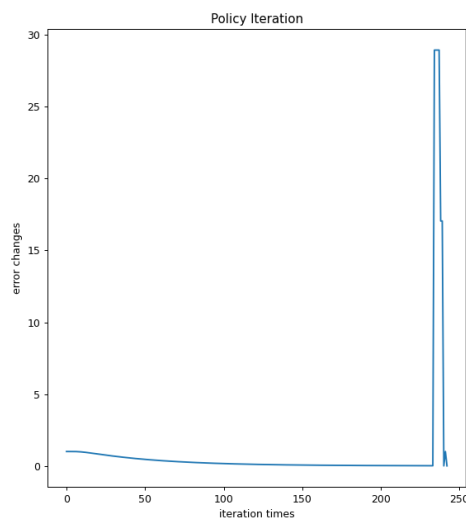
- Value Iteration

-1.0	0.0	-1.0	-2.0	-3.0	-4.0
-2.0	-1.0	-2.0	-3.0	-4.0	-4.0
-3.0	-2.0	-3.0	-4.0	-4.0	-3.0
-4.0	-3.0	-4.0	-4.0	-3.0	-2.0
-5.0	-4.0	-4.0	-3.0	-2.0	-1.0
-5.0	-4.0	-3.0	-2.0	-1.0	0.0

→	terminal	←	←	←	←
↖	↑	↖	↖	↖	↓
↖	↑	↖	↖	↘	↓
↖	↑	↖	↘	↘	↓
↖	↑	↘	↘	↘	↓
→	→	→	→	→	terminal

From above figures, we can see that we get the correct values and actions for all states.

What's more, we can compare the performance of the two methods.



```
PS E:\大三下\强化学习\作业\homework1> D:\python\python.exe "e:\大三下\强化学习\作业\homework1\solution.py"
The spent time of Policy Iteration: 0.05807638168334961
The spent time of Value Iteration: 0.00799703598022461
```

- From the first figure, we can see that the iteration times of value iteration is much smaller than the policy iteration. And the error change of policy iteration is much bigger than the value iteration.
- From the second figure, we can see that the spent time of policy iteration is almost ten times bigger than value iteration's.
- From above results, we can conclude that the performance of value iteration is much better than policy iteration in this problem. However, this doesn't mean value iteration is always better than policy iteration. This depends on specific problems.