

Report for Gridworld

516021910212 Liuyuxi

March 2019

1 Introduction

Use dynamic programming to carry out reinforcement learning. Judge the value function of a strategy based on known models. And find the optimal strategy and optimal value function by policy evaluation and policy iteration. The DP methods have been tested in a small gridworld.

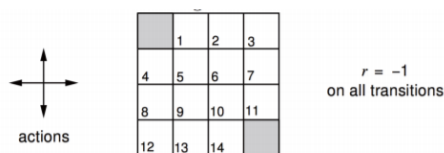


Figure 1: The gridworld

The agent is in an 4×4 gridworld, trying to get to the top left or bottom right corner as shown in 1. Each grid in the gridworld represents a certain state. S_t is the state at grid t . And the state space can be denoted as $S = \{s_t | t \in 0, \dots, 15\}$. S_0 and S_{15} are terminal states, where S_1 to S_{14} are non-terminal states and can move one grid to north, east, south and west. Hence the action space is $A = \{N, E, S, W\}$. And the actions leading out of the grid leave state unchanged. Each movement get a reward of -1 until the terminal state is reached. Both policy evaluation and policy iteration have been used in finding the shortest way to the terminal state randomly given an initial non-terminal state.

2 Algorithms

This is a model in the Reinforcement Learning literature. In this particular case:

- State space S : GridWorld has 4×4 distinct states, The start state is S_0 and the end state is S_{15}

- Actions A: The agent can choose from up to 4 actions to move around.
- Transition probability P: The position of any attempt to leave the square world will not change, and the rest will be 100% transferred to the position pointed by the action
- Rewards R: The agent receives -1 reward when it is in the non-terminal states and receives 0 reward when it is the terminal states.
- Attenuation coefficient $\gamma : 1$
- Current strategy : The individual adopts a random action strategy and has an equal probability to move in any possible direction in any non-terminating state, $\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 1/4$

In other words, this is a deterministic, finite Markov Decision Process (MDP) and as always the goal is to find an agent policy that maximizes the future discounted reward.

- The Policy Evaluation (one sweep) button turns the Bellman equation into a synchronous update and performs a single step of Value Function estimation. Intuitively, this update is diffusing the raw Rewards at each state to other nearby states through 1. the the dynamics of the environment and 2. the current policy.

$$V^\pi(s) = \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

- The Policy Update button iterates over all states and updates the policy at each state to take the action that leads to the state with the best Value (integrating over the next state distribution of the environment for each action).

Here, we are computing the action value function at each state $Q(s, a)$, which measures how much expected reward we would get by taking the action a in the state s . The function has the form:

$$Q^\pi(s, a) = E_\pi\{r_t + \gamma V^\pi(s_{t+1}) \mid s_t = s, a_t = a\} = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

In our Gridworld example, we are looping over all states and evaluating the Q function for each of the (up to) four possible actions. Then we update the policy to take the argmaxy actions at each state. That is,

$$\pi'(s) = \arg \max_a Q(s, a)$$

3 Setup

3.1 Policy Evaluation

```

1  #declare the state
2  states=[i for i in range(16)]
3  #declare the value and the initialization is 0
4  values=[0 for _ in range(16) ]
5  actions=["N","E","S","W"]
6  #move north then the current state -4
7  ds_actions={"N":-4,"E":1,"S":4,"W":-1}
8  gamma=1.00#Attenuation coefficient
9  #determin the next state
10 def nextState(s,a):
11     next_state=s
12     #boundary
13     if (s%4==0 and a=="W") or (s<4 and a=="N") \
14     or ((s+1)%4==0 and a=="E") or (s>11 and a=="S"):
15         pass
16     else:
17         ds=ds_actions[a]
18         next_state=s+ds
19     return next_state
20
21 #get the reward when leaving the state
22 def rewardOf(s):
23     return 0 if s in [0,15] else -1
24
25 #check whether the terminate state
26 def isTerminateState(s):
27     return s in [0,15]
28
29 #get all the next states
30 def getAllState(s):
31     next=[]
32     if isTerminateState(s):
33         return next
34     for a in actions:
35         next_state=nextState(s,a)
36         next.append(next_state)#include not move
37     return next
38
39 #update the current value
40 def updateValue(s):
41     AllState=getAllState(s)
42     newValue=0
43     num=4
44     reward=rewardOf(s)#leaving reward
45     for next_state in AllState:

```

```

46         newValue+=1.00/num*(reward+\
47         gamma*values[next_state])
48     return newValue
49
50
51 def printValue(v):
52     for i in range(16):
53         print('{0:>6.2f}'.format(v[i]), end=" ")
54         if (i + 1) % 4 == 0:
55             print("")
56     print()
57
58
59 def performOneIteration():
60     newValues=[0 for _ in range(16)]
61     for s in states:
62         newValues[s]=updateValue(s)
63     global values
64     values=newValues
65     printValue(values)
66
67 def main():
68     max_iterate_times=160
69     cur_iterate_times=0
70     while cur_iterate_times<=max_iterate_times:
71         print("Iterate_No.{0}".format\
72         (cur_iterate_times))
73         performOneIteration()
74         cur_iterate_times+=1
75     printValue(values)
76
77 if __name__=='__main__':
78     main()

```

3.2 Policy Iteration

```

1  # the state
2  states = [i for i in range(16)]
3  # the value
4  values = [0 for _ in range(16)]
5  actions = ["N", "E", "S", "W"]
6  # move north then the current state -4
7  ds_actions = {"N": -4, "E": 1, "S": 4, "W": -1}
8  gamma = 1.00 # Attenuation coefficient

```

```

9  a = []
10 for i in range(16):
11     a.append(actions)
12 #determin the next state
13 def nextState(s,a):
14     next_state=s
15     #boundary
16     if (s%4==0 and a=="W") or (s<4 and a=="N") or \
17     ((s+1)%4==0 and a=="E") or (s>11 and a=="S"):
18         pass
19     else:
20         ds=ds_actions[a]
21         next_state=s+ds
22     return next_state
23
24 #get the reward when leaving the state
25 def rewardOf(s):
26     return 0 if s in [0,15] else -1
27
28
29
30 #check whether the terminatestate
31 def isTerminateState(s):
32     return s in [0,15]
33
34
35 def printValue(v):
36     for i in range(16):
37         print( '{0:>6.2f}'.format(v[i]),end = " ")
38         if (i+1)%4==0:
39             print("")
40     print()
41
42
43 #strategy
44 def allowedActions(s):
45     allow=[]
46     candidate=[]
47     if isTerminateState(s):
48         return allow
49     for a in actions:
50         next_state=nextState(s,a)
51         candidate.append(values[next_state])
52     b=max(candidate)
53     for a in actions:
54         next_state=nextState(s,a)

```

```

55         if(values[next_state]==b):
56             allow.append(a)
57     return allow
58
59 #get the possible state on allowed actions
60 def allowedStates(s):
61     allowStates=[]
62     if isTerminateState(s):
63         return allowStates
64     for a in allowedActions(s):
65         next_state=nextState(s,a)
66         allowStates.append(next_state)
67     return allowStates
68
69 def allowedUpdateValue(s):
70     next_state=allowedStates(s)
71     newValue=0
72     num=len(next_state)
73     reward=rewardOf(s)
74     for state in next_state:
75         newValue+=1.00/num*(reward+\
76             gamma*values[state])
77     return newValue
78
79 #iteration and update the strategy
80 def allowedPerformOneIteration():
81     newValues=[0 for _ in range(16)]
82     newa=[]
83     for s in states:
84         allowedActions(s)
85         newa.append(allowedActions(s))
86     global a
87     a=newa
88     #the strategy is updated
89     for s in states:
90         newValues[s]=allowedUpdateValue(s)
91     global values
92     values=newValues
93     printValue(values)
94
95 def main():
96     max_iterate_times=10
97     cur_iterate_times=0
98     while cur_iterate_times<=max_iterate_times:
99         print("Iterate_No.{0}".format\
100             (cur_iterate_times))

```

```

101         allowedPerformOneIteration()
102         #performOneIteration()
103         cur_iterate_times+=1
104         printValue(values)
105         for i in range(16):
106             print(allowedActions(i))
107
108
109 if __name__=='__main__':
110     main()

```

4 Result

4.1 Policy Evaluation

With the same action strategy, the algorithm will converge around 153 times. We set the maximum number of iterations to 160. The resulting value function is shown in Table1 and Figure2.

Table 1: Policy Evaluation

Iterate No.152			
0.00	-14.00	-20.00	-21.99
-14.00	-18.00	-20.00	-20.00
-20.00	-20.00	-18.00	-14.00
-21.99	-20.00	-14.00	0.00
Iterate No.153			
0.00	-14.00	-20.00	-22.00
-14.00	-18.00	-20.00	-20.00
-20.00	-20.00	-18.00	-14.00
-22.00	-20.00	-14.00	0.00
Iterate No.154			
0.00	-14.00	-20.00	-22.00
-14.00	-18.00	-20.00	-20.00
-20.00	-20.00	-18.00	-14.00
-22.00	-20.00	-14.00	0.00

4.2 Policy Iteration

the algorithm will converge around 4 times. We set the maximum number of iterations to 10. The result is shown in Table2 and Figure3.

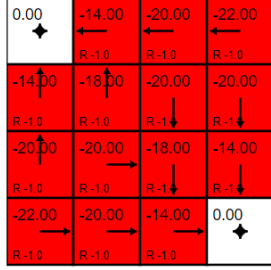


Figure 2: Policy evaluation

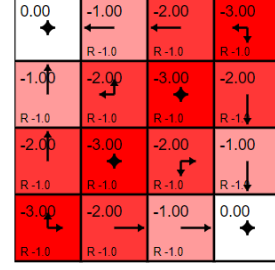


Figure 3: Policy Iteration

Table 2: Policy iteration

Iterate No.2			
0.00	-1.00	-2.00	-2.00
-1.00	-2.00	-2.00	-2.00
-2.00	-2.00	-2.00	-1.00
-2.00	-2.00	-1.00	0.00
Iterate No.3			
0.00	-1.00	-2.00	-3.00
-1.00	-2.00	-3.00	-2.00
-2.00	-3.00	-2.00	-1.00
-3.00	-2.00	-1.00	0.00
Iterate No.4			
0.00	-1.00	-2.00	-3.00
-1.00	-2.00	-3.00	-2.00
-2.00	-3.00	-2.00	-1.00
-3.00	-2.00	-1.00	0.00

5 Conclusion

Thanks to this project that I have a deep understanding of reinforcement learning. I also meet with some obstacles during the period. At the beginning, I choose the environment gridworld. However, after some investigation I turn to establish the environment myself. Thanks for the guidance of the teacher and our teaching assistants.