

本科实验报告

课程名称: 人工智能

实验名称: 值迭代与策略迭代

专业名称: 保密管理

学生姓名: 武自厚

学生学号: 20336014

实验地点: 东校园实验楼 D502

实验成绩:

报告时间: 2022 年 6 月 14 日

一、 实验要求

基于 gym 库中的 Frozen Lake 环境在 vi_and_pi.py 中实现策略迭代和值迭代算法,并输出算法 收敛后的路径.

- 实现 vi_and_pi.py 中的策略评估、策略提升、策略迭代函数.
- 实现 vi_and_pi.py 中的值迭代函数.
- 需要提交一份简要报告 + 代码.

二、 实验过程

1. 强化学习原理

强化学习,是机器学习的范式和方法论之一. 用于描述和解决智能体 (agent) 在与环境的交互过程中通过学习策略以达成回报最大化或实现特定目标的问题. 强化学习的常见模型是标准的 Markov 决策过程.Markov 决策过程可以用六元组 (S,A,R,T,P_0,γ) 描述,其中

- S表示状态空间。
- A 表示行为空间.
- R = R(s, a) 表示在特定状态下进行特定行为的奖励值函数.
- $T: S \times A \times S \mapsto [0,1]$ 表示状态变化函数.
- P_0 表示初始状态的分布.
- γ表示折扣因子.
- 目标是找到一个能使奖励值期望最大化的策略.

2. 策略迭代原理

从一个初始化的策略出发, 先进行策略评估, 然后改进策略, 评估改进的策略, 再进一步改进策略. 经过不断迭代更新, 直达策略收敛, 这种算法被称为"策略迭代".

此算法可以用伪代码表示:

Algorithm 1: 策略评估算法 evaluate()

输入: Markov 六元组 $(S, A, R, T, P_0, \gamma)$, 临界误差 θ

输出: 评估价值映射 V 1 从 P_0 中获取价值映射 V

2 do

```
\begin{array}{c|cccc} \mathbf{3} & \Delta := 0 \\ \mathbf{4} & \mathbf{for} \ s \in S \ \mathbf{do} \\ \mathbf{5} & v := V(s) \\ \mathbf{6} & V(s) := R(s, \pi(s)) + \gamma \sum_{s'} T(s, a, s') V(s') \\ \mathbf{7} & \Delta := \max\{\Delta, |v - V(s)|\} \end{array}
```

- 8 until $\Delta < \theta$;
- 9 return V

Algorithm 2: 策略改进算法 improve()

输入: Markov 六元组 $(S, A, R, T, P_0, \gamma)$, 现有策略 π , 已有价值 V

输出: 策略迭代结果 π'

1 for $s \in S$ do

2 $\pi(s) := \arg \max_{a \in A} \left\{ R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s') \right\}$

з return π

Algorithm 3: 策略迭代算法

输入: Markov 六元组 $M = (S, A, R, T, P_0, \gamma)$, 临界误差 θ

输出: 价值评估结果 V, 策略迭代结果 π

1 随机初始化策略 π

2 do

```
\begin{array}{c|cccc} \mathbf{3} & V := \mathtt{evaluate}(M, \theta) \\ \mathbf{4} & \pi := \mathtt{improve}(M, \pi, V) \end{array}
```

- 5 until π 收敛;
- 6 return V, π

3. 价值迭代原理

对每一个当前状态 s, 对每个可能的动作 a 都计算一下采取这个动作后到达的下一个状态的期望价值. 看看哪个动作可以到达的状态的期望价值函数最大, 就将这个最大的期望价值函数作为当前状态的价值函数 V(s), 循环执行这个步骤, 直到价值函数收敛.

此算法可以用伪代码表示:

Algorithm 4: 价值迭代算法

```
输入: Markov 六元组 M = (S, A, R, T, P_0, \gamma), 临界误差 \theta
   输出: 价值评估结果 V, 策略迭代结果 \pi
 1 随机初始化策略 \pi, 零初始化价值函数 V
 2 do
       V' := V
 3
       \Delta := 0
 4
       for s \in S do
           Q_{\pi}(a) := R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s')
 6
           V(s) := \max_{a \in A} Q_{\pi}(a)
 7
           \pi(s) := \arg \max_{a \in A} Q_{\pi}(a)
           \Delta := \max\{\Delta, |V(s) - V'(s)|\}
10 until \Delta < \theta;
11 return V, \pi
```

4. 代码实现

```
def policy_evaluation(P, nS, nA, policy, gamma=0.9, tol=1e-3):
        """Evaluate the value function from a given policy.
44
45
       Parameters
46
        _____
47
       P, nS, nA, gamma:
48
            defined at beginning of file
49
       policy: np.array[nS]
            The policy to evaluate. Maps states to actions.
        tol: float
52
            Terminate policy evaluation when
53
                max |value_function(s) - prev_value_function(s)| < tol</pre>
        Returns
55
56
        value_function: np.ndarray[nS]
            The value function of the given policy, where value function[s] is
            the value of state s
59
        11 11 11
60
       value_function = np.zeros(nS)
62
        i = 0
63
        while True:
64
            prev_value_function = np.copy(value_function)
            for state in range(nS):
66
                action = policy[state]
67
```

姓名:武自厚

```
value_function[state] = sum(
68
                      [prob * (reward + gamma * prev_value_function[next_state])
69
                      for prob, next_state, reward, done in P[state][action]]
70
71
             if np.max(np.fabs(value_function - prev_value_function)) < tol:</pre>
72
                 break
73
             else:
74
                 i += 1
75
76
        return value_function
77
78
79
    def policy_improvement(P, nS, nA, value_from_policy, policy, gamma=0.9):
80
         """Given the value function from policy improve the policy.
82
        Parameters
83
         _____
84
        P, nS, nA, gamma:
             defined at beginning of file
86
        value_from_policy: np.ndarray
87
             The value calculated from the policy
88
        policy: np.array
89
             The previous policy.
90
91
        Returns
         -----
93
        new_policy: np.ndarray[nS]
94
             An array of integers. Each integer is the optimal action to take
95
             in that state according to the environment dynamics and the
96
             given value function.
97
         11 11 11
98
        new_policy = np.zeros(nS, dtype='int')
100
101
        for state in range(nS):
102
             q_value = np.zeros(nA)
103
             for action in range(nA):
104
                 for next_sr in P[state][action]:
105
                     prob, next_state, reward, done = next_sr
106
                     q_value[action] += prob * (reward + gamma *
107

    value_from_policy[next_state])

             new_policy[state] = np.argmax(q_value)
108
109
```

```
return new_policy
110
111
112
    def policy_iteration(P, nS, nA, gamma=0.9, tol=10e-3):
113
         """Runs policy iteration.
114
115
         You should call the policy_evaluation() and policy_improvement() methods to
116
         implement this method.
117
118
        Parameters
119
         _____
120
        P, nS, nA, gamma:
121
             defined at beginning of file
122
         tol: float
123
             tol parameter used in policy_evaluation()
124
        Returns:
125
         _____
126
        value_function: np.ndarray[nS]
127
        policy: np.ndarray[nS]
128
         n n n
129
130
        policy = np.zeros(nS, dtype=int)
131
        i = 0
132
        while True:
133
             value_function = policy_evaluation(P, nS, nA, policy)
             next_policy = policy_improvement(P, nS, nA, value_function, policy)
135
             if np.all(policy == next_policy):
136
                 print(f'Policy Iteration converged as step {i}')
137
                 break
138
             else:
139
                 policy = next_policy
140
                 i += 1
        return value_function, policy
142
143
144
    def value_iteration(P, nS, nA, gamma=0.9, tol=1e-3):
145
146
        Learn value function and policy by using value iteration method for a given
147
        gamma and environment.
148
149
        Parameters:
150
         _____
151
         `P`, `nS`, `nA`, `gamma`:
152
```

```
defined at beginning of file
153
         tol: float
154
             Terminate value iteration when
155
                 max |value_function(s) - prev_value_function(s)| < tol</pre>
156
        Returns:
157
         _____
158
        value_function: np.ndarray[nS]
159
        policy: np.ndarray[nS]
160
         n n n
161
162
        value_function = np.zeros(nS)
163
        policy = np.zeros(nS, dtype=int)
164
        i = 0
165
        while True:
             prev_value_function = value_function.copy()
167
168
             for state in range(nS):
169
                 q_values = np.zeros(nA)
                 for action in range(nA):
171
                     for next_sr in P[state][action]:
172
                          prob, next_state, reward, done = next_sr
173
                          q_values[action] += prob * (reward + gamma *
174
                          → prev_value_function[next_state])
                 value_function[state] = np.max(q_values)
175
                 policy[state] = np.argmax(q_values)
177
             i += 1
178
179
             if np.max(np.fabs(value_function - prev_value_function)) < tol:</pre>
180
                 print(f"Convergence at iteration {i}")
181
                 break
182
183
        return value_function, policy
184
```

三、 实验结果

在 Frozen Lake 问题应用价值迭代和策略迭代得到结果:

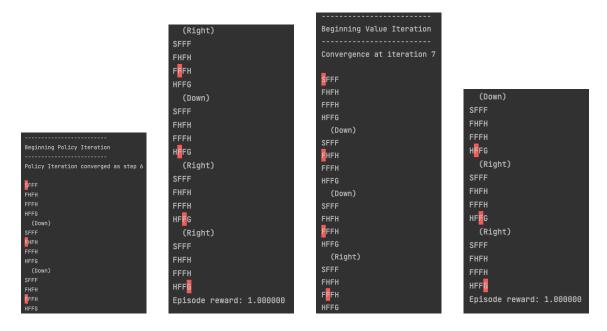


图 1: Frozen Lake 问题中的强化学习 agent 结果

可以看到 agent 在价值迭代和策略迭代中都找出了最好的策略. 本次实验成功.