



گزارش فاز اول پروژه
مبانی و کاربردهای هوش مصنوعی
دکتر حسین کارشناس

آدرینا ابراهیمی ۹۹۳۶۲۳۰۰۲

کیان مجلسی ۹۹۳۶۱۳۰۵۱

اردیبهشت ۱۴۰۲

نام و علت انتخاب الگوریتم

از الگوریتم تکرار ارزش (Value iteration) استفاده شده است که مطابق با شبه کد زیر بر اساس سودمندی مورد انتظار سیاست بهینه را پیدا می‌کند. به علت سادگی در پیاده‌سازی و همگرا شدن به یک جواب بهینه از این الگوریتم استفاده شده است.

```
 $V^*, Q^* = \text{VALUE-ITERATION}(\text{MDP})$   
// MDP consist of States, Actions(s),  $R(s, a, s')$ ,  $T(s'|s, a)$ , and discount factor  $\gamma$   
hyperparameters: maxIters //maximum number of steps in the environment  
  
foreach state  $s \in \text{States}$  do  
     $V^*(s) \leftarrow 0$   
    foreach action  $a \in \text{Actions}(s)$  do  
         $Q^*(s, a) \leftarrow 0$   
  
     $t = 1$   
    while  $t < \text{maxIters}$  and not (converged or terminated) do  
        foreach state  $s \in \text{States}$  do // except terminal states  
            foreach action  $a \in \text{Actions}(s)$  do // old values buffered  
                 $Q^*(s, a) \leftarrow \sum_{s'} T(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$   
                 $V^*(s) \leftarrow \max_a Q^*(s, a)$   
  
         $t = t + 1$   
return  $V^*, Q^*$ 
```

45

گزارش کار الگوریتم

تابع `value_iteration`:

در این تابع عملیات اصلی الگوریتم ذکر شده انجام شده است. در خطوط ۷ تا ۱۶ دیکشنری‌های `v_value` (مقدار مورد انتظار سودمندی با شروع از هر حالت) و `q_value` (مقدار مورد انتظار سودمندی با انجام هر کنش در هر حالت) مقداردهی اولیه شده‌اند و خانه‌های چاله مشخص شده‌اند. در خط ۱۹ مقدار مورد انتظار سودمندی برای خانه هدف برای ۱ در نظر گرفته شده و در خط ۲۰ و ۲۱ این مقدار برای خانه‌های چاله برابر ۱- در نظر گرفته شده است.

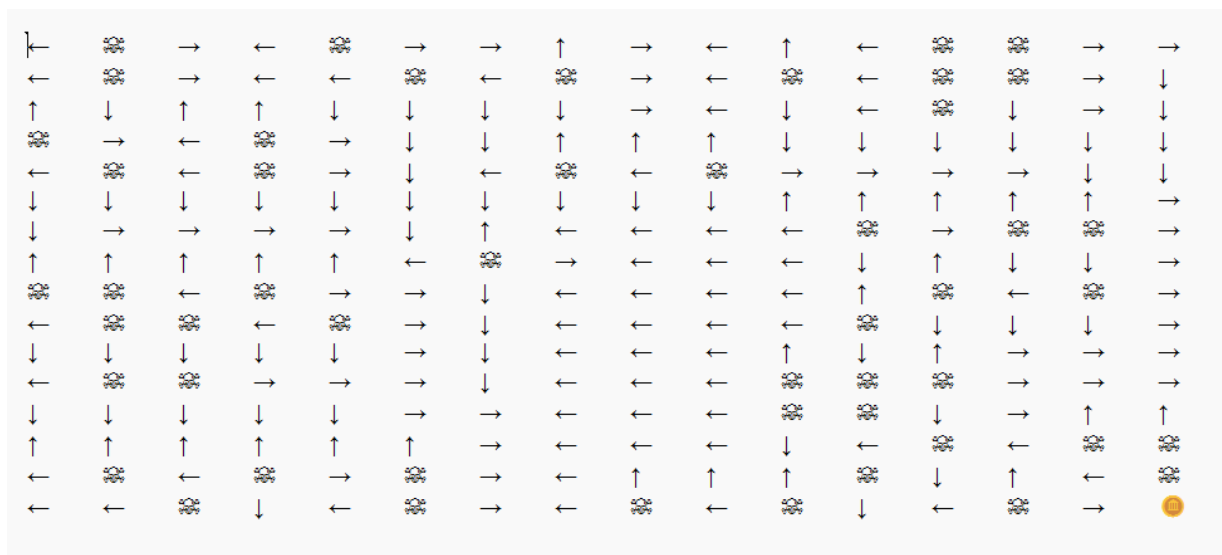
سپس در یک حلقه به طول ۱۰۰۰۰ (برای همگرایی سیاست‌های ممکن به یک حالت بهینه)، یک حلقه روی تمام حالت‌های موجود در محیط زده و پس از آن، برای تمام کنش‌های موجود در آن حالت احتمال، پاداش و سودمندی مورد انتظار با شروع از حالت را برای نتیجه‌های ممکن آن کنش بررسی کرده و q_value را برای آن کنش در آن حالت محاسبه می‌کنیم. لازم به ذکر است پاداش خانه‌های چاله برابر ۱-، پاداش خانه‌ی هدف برابر ۱ و پاداش سایر خانه‌ها با توجه به فاصله‌شان تا خانه هدف و تابع لجستیک مقداری بین ۰ و ۱- در نظر گرفته شده است. در نهایت مقدار مورد انتظار سودمندی با شروع از آن حالت (v_value) و سیاستی که باید در پیش گرفته شود را آپدیت کرده و در صورت برقرار شدن شرط همگرایی از حلقه‌ها خارج می‌شویم.

```

1  def value_iteration():
2      # Initialize
3      v_values, q_values = {}, {}
4      convergenceTrack = [0]
5
6      # Initialize the value of each state to 0 and store the terminal states
7      for state in env.P:
8          if state == goal_state:
9              continue
10         v_values[state] = 0
11         q_values[state] = {}
12         for act in env.P[state]:
13             q_values[state][act] = 0
14             for probability, nextState, reward, isTerminalState in env.P[state][act]:
15                 if (reward == 0) and isTerminalState:
16                     terminal_states.add(nextState)
17
18     # Set the value of the goal state to 1 and the terminal states to -1
19     v_values[goal_state] = 1
20     for ts in terminal_states:
21         v_values[ts] = -1
22
23     for i in range(10000):
24         # Check states in the environment
25         for state in env.P:
26             if (state not in terminal_states) and (state != goal_state):
27                 # Check actions in that state
28                 for act in env.P[state]:
29                     s = 0
30                     # Calculate every result of the action
31                     for probability, nextState, reward, isTerminalState in env.P[state][act]:
32                         # Calculate the reward of each actions
33                         if (reward == 0) and isTerminalState:
34                             reward = -1
35                         elif (reward == 1) and isTerminalState:
36                             reward = 1
37                         else:
38                             # Calculate the distance to the goal
39                             x, y = nextState // map_size, nextState % map_size
40                             dist_to_goal = np.sqrt(
41                                 np.power(x - (map_size - 1), 2) + np.power(y - (map_size - 1), 2))
42                             reward = -0.1 / (1 + np.exp(-dist_to_goal))
43                     s += probability * \
44                         (reward + (discount_factor * v_values[nextState]))
45                 # Update the q_value of state and action
46                 q_values[state][act] = s
47             # Update the v_value of state and policy
48             v_values[state] = max(q_values[state].values())
49             # Check convergence
50             convergenceTrack.append(
51                 np.linalg.norm(List(v_values.values())))
52             if (i > 1000) and np.isclose(convergenceTrack[-1], convergenceTrack[-2]):
53                 print('Values Converged')
54                 return v_values, q_values
55     return v_values, q_values

```

پس از محاسبه سیاست بهینه یک فایل که حاوی نقشه سیاست بهینه است ایجاد می‌شود که محتویات آن به شکل زیر می‌باشد.



منابع ایده

از منابع زیر برای آشنایی با نحوه پیاده‌سازی و الگو گرفتن اولیه استفاده شده است.

Stuart J. Russell, Peter Norvig - Artificial Intelligence_ A Modern Approach, Global Edition-Pearson (2021)

[Frozen Lake - Gymnasium Documentation \(farama.org\)](https://farama.org/)