



گزارش فاز دوم پروژه: دریاچه یخی

مبانی و کاربردهای هوش مصنوعی

دکتر حسین کارشناس

آدرینا ابراهیمی ۹۹۳۶۲۳۰۰۲

کیان مجلسی ۹۹۳۶۱۳۰۵۱

اردیبهشت ۱۴۰۲

## نام و علت انتخاب الگوریتم

از الگوریتم Q-Learning که Off-policy (مقادیر جدول  $q\_table$  را متفاوت از سیاست رفتاری  $\epsilon$ -greedy آپدیت می‌کند.) است استفاده شده زیرا نسبت به الگوریتم SARSA که On-policy (مقادیر جدول  $q\_table$  را مشابه سیاست رفتاری  $\epsilon$ -greedy آپدیت می‌کند.) است تخمین بهتری از سیاست بهینه به دست می‌آورده، مستقیماً از سیاست بهینه یاد می‌گیرد و تصمیم‌گیری بهینه‌تری نسبت به الگوریتم SARSA دارد؛ این در حالی است که SARSA از سیاستی نزدیک به سیاست بهینه یاد می‌گیرد و برای محاسبه پاداش حالت‌های بعدی، نیازی به داشتن تمام  $q\_table$  نیست

## گزارش کار الگوریتم

### تابع $Q\_Learning\_Algorithm$ :

در این تابع عملیات اصلی الگوریتم Q Learning انجام می‌شود. این تابع جدول Q table و مشاهدات از محیط را دریافت می‌کند. سپس، در یک حلقه به طول حداکثر تعداد دفعاتی که می‌خواهیم عمل یادگیری الگوریتم انجام شود، یک حلقه دیگر به طول حداکثر تعداد حرکات در محیط می‌زنیم. در حلقه دوم با توجه به مشاهداتی که از محیط داریم با استفاده از سیاست  $\epsilon$ -greedy یک کنش انتخاب می‌کنیم. پس از آن، با استفاده از تابع  $step$  کنش انتخاب شده را در محیط اعمال کرده و این تابع مشاهده محیط، پاداش، خاتمه یافتن یا نیافتن، متوقف شدن یا نشدن و اطلاعات محیط را برمی‌گرداند.

مقدار مورد انتظار سودمندی در صورت افتادن آدمک در چاله ۱- و در صورت رسیدن آدمک به مقصد ۱ در نظر گرفته شده است. در صورتی که آدمک روی خانه‌های یخ‌زده باشد مقدار مورد انتظار سودمندی برابر فاصله اقلیدسی مختصات آدمک تا مختصات هدف است.

سپس، مقدار Q table را برای آن حالت و کنش در محیط با توجه به فرمول آپدیت مقدار Q table در الگوریتم Q Learning آپدیت می‌کنیم. در صورت پایان یا متوقف شده حرکت آدمک در محیط از حلقه دوم خارج شده و محیط را ریست می‌کنیم. در نهایت شرط همگرایی را بررسی که و در صورتی که مقادیر جدول Q Learning تغییر چندانی نداشته باشند از حلقه اول نیز  $for$  نیز خارج می‌شویم.

```

1 def Q_Learning_Algorithm(observation, q_table):
2     # train the agent for max_iter_number number of episodes
3     for i in tqdm(range(0, max_iter_number)):
4
5         # reset the environment
6         observation, info = env.reset()
7         terminated, truncated = False, False
8
9         # for each episode, run the algorithm for MAX_STEP number of steps
10        for __ in range(MAX_STEP):
11            # choose the action using epsilon greedy policy
12            action = epsilon_greedy_policy(q_table, observation)
13
14            next_observation, reward, terminated, truncated, info = env.step(action)
15
16            # if the agent reaches the goal, then reward = 1, else if it falls to a hole, reward = -1
17            if (reward == 0) and terminated:
18                reward = -1
19            elif (reward == 1) and terminated:
20                reward = 1
21            else:
22                # Calculate the distance to the goal for agent and reward it based on the distance
23                x, y = next_observation // map_size, next_observation % map_size
24                dist_to_goal = np.sqrt(np.power(x - (map_size - 1), 2) + np.power(y - (map_size - 1), 2))
25                reward = -0.1 / (1 + np.exp(-dist_to_goal))
26
27            # choose the next action
28            next_action = np.argmax(q_table[next_observation])
29
30            # update the q_table
31            q_table[observation][action] = q_table[observation][action] + LEARNING_RATE * \
32                (reward + DISCOUNT_FACTOR * q_table[next_observation][next_action] - q_table[observation][action])
33
34            # update the observation
35            observation = next_observation
36
37            # if the agent reaches the goal or falls into a hole, then reset the environment
38            if terminated or truncated:
39                observation, info = env.reset()
40                break
41
42            convergenceTrack.append(np.linalg.norm(q_table.flatten().tolist()))
43            if (i >= 1000) and np.isclose(convergenceTrack[-1], convergenceTrack[-2]):
44                print('\nValues Converged')
45                return
46        print("Training Completed")
47        sleep(2)

```

## تابع epsilon\_greedy\_policy:

در این تابع یک عدد رندوم تولید شده و در صورتی که آن عدد کوچکتر از اپسیلون باشد، به صورت رندوم یک کنش انتخاب کرده و در صورتی که بزرگتر از اپسیلون باشد بهترین کنش ممکن را از جدول Q table انتخاب می‌کنیم.



```

1 def epsilon_greedy_policy(q_table, observation):
2     if np.random.uniform(0, 1) < EPSILON:
3         action = env.action_space.sample()
4     else:
5         action = np.argmax(q_table[observation])
6     return action

```

## تابع save\_q\_table:

در این تابع ابتدا حالت‌های پایانی محیط را به دست آورده سپس، با استفاده از جدول Q table و با یک حلقه روی تمام حالت‌های محیط نقشه محیط به همراه policy به دست آورده شده را رسم می‌کنیم.

```
1 def save_q_table(q_table):
2     terminal_states = set()
3     goal_state = (map_size - 1) * map_size + (map_size - 1)
4
5     for state in env.P:
6         if state == goal_state:
7             continue
8         for act in env.P[state]:
9             for probability, nextState, reward, isTerminalState in env.P[state][act]:
10                 if (reward == 0) and isTerminalState:
11                     terminal_states.add(nextState)
12
13     with open('q_table_frozen.txt', 'w', encoding="utf-8") as inp:
14         for state in range(map_size**2):
15             if state in terminal_states:
16                 inp.write(u'❌\t')
17             elif state == goal_state:
18                 inp.write(u'🏁\t')
19
20             else:
21                 if np.all(q_table[state] == 0):
22                     inp.write(u'□\t')
23                 else:
24                     argm = np.argmax(q_table[state])
25                     if argm == 0:
26                         inp.write(u'←\t')
27                     elif argm == 1:
28                         inp.write(u'↓\t')
29                     elif argm == 2:
30                         inp.write(u'→\t')
31                     elif argm == 3:
32                         inp.write(u'↑\t')
33             if (state + 1) % map_size == 0:
34                 inp.write('\n')
```

## قسمت main برنامه:

در ابتدا قسمت train را با استفاده الگوریتم Q Learning انجام شده و سیاست به دست آمده را ذخیره می‌کنیم.

پس از آن محیط را ریست کرده و در یک اپیزود سیاست به دست آمده را تست می‌کنیم. و در نهایت پس از پایان اپیزود نمودار همگرایی را رسم می‌کنیم.



```

1  if __name__ == "__main__":
2
3      # train the agent with no graphical output
4      Q_Learning_Algorithm(observation, q_table)
5
6      save_q_table(q_table)
7
8      env.close()
9
10     # create Enviroment
11     env = gym.make("FrozenLake-v1", desc=map_shape, render_mode="human", is_slippery=False)
12
13     # reset enviroment
14     observation, info = env.reset()
15
16     # test the agent for max_iter_number number of episodes
17     for i in range(1):
18
19         terminated, truncated = False, False
20
21         for _ in range(MAX_STEP):
22             # select the action
23             action = np.argmax(q_table[observation])
24
25             next_observation, reward, terminated, truncated, info = env.step(action)
26
27             env.render()
28
29             # check if the agent reaches the goal or falls into a hole
30             if terminated or truncated:
31                 observation, info = env.reset()
32                 break
33
34             observation = next_observation
35
36
37     env.close()
38     plot_convergence(convergenceTrack)

```

## تحلیل نتایج

با تغییر مقدار اپسیلون از  $\frac{0}{9}$  به سمت  $\frac{0}{2}$  مشاهده می‌شود سیاست به دست آمده توسط الگوریتم بهبود پیدا کرده و به سمت سیاست بهینه تغییر پیدا می‌کند. شکل اول مقدار اپسیلون برابر  $\frac{0}{9}$  و شکل دوم مقدار اپسیلون برابر  $\frac{0}{5}$  و شکل سوم مقدار اپسیلون برابر  $\frac{0}{2}$  است. خانه‌های سفید نشان‌دهنده این است که عامل تا به حال در این خانه‌ها نرفته و خانه‌هایی که اسکلت دارند نشان‌دهنده چاله و خانه سکه نشان‌دهنده خانه هدف هستند. همانطور که مشاهده می‌شود با کاهش مقدار اپسیلون خانه‌های بیشتری کاوش شده و سیاست تکمیل می‌شود.

☞ □ □ □ □ ☞ □ □ ☞ ☞ ☞ ☞ ☞ □ □ ☞

□ □ ☞ □ □ ☞ □ □ ☞ ☞ □ □ □ □

□ □ □ □ □ □ □ □ □ □ ☞ ☞ □ ☞

□ □ □ □ ☞ ☞ ☞ □ □ □ □ □ □

□ □ □ ☞ □ ☞ □ ↓ → □ ☞ □ □ □ □

□ ↑ → ↓ ☞ □ □ ↓ → → ↓ ☞ □ ☞ ☞ □

↓ → ← ↓ ☞ □ ☞ ↓ ↓ → ↓ □ ↓ ↓ → □

↑ ← ☞ □ □ ↑ ☞ □ ☞ ← ↓ → ↓ ↓ □ □

↑ → → ↑ ☞ ↑ → ↑ ↓ ↓ ↓ ↓ ↓ □ ☞

← → → ↑ ← ↑ → ↑ ↓ ☞ ☞ □ ☞ □ ☞

☞ ↑ ↑ ☞ ↑ ↓ → ↓ ☞ ↓ ↓ ☞ □ ☞ □ ☞

→ ↑ ↑ → → → → → ↓ □ □ □ ☞ □

↑ ☞ ↑ ↑ ↑ ↑ ↑ ☞ ↓ → □ ↑ ↓ □ □ ☞

↑ ← ☞ → → → → → → ☞ ☞ → □ □ □

↑ → → → → → ↓ ☞ ↑ ↓ ↓ ← → □ ☞ ☞

↑ → → → ↑ ☞ → → ↑ ← ↓ ↓ → ☞ □ □

☞ □ □ □ □ ☞ ☞ ☞ □ ☞ ☞ □ □ □ ☞

□ ↑ ↑ → ↑ → ☞ □ □ □ ☞ □ □ □ □

↓ ← → ↑ ↓ ☞ → □ □ □ ☞ □ □ □ □

☞ → → ↑ ← → → ↑ ☞ □ □ □ □ □ □

☞ ↓ → → → ↓ → → → □ □ □ □ ☞ □

☞ ↑ ↑ ↑ → → ← → ☞ ☞ □ □ ☞ □ □ □

← ← ☞ → ↑ → ↑ ☞ □ ☞ □ ☞ ☞ □ □ □

↑ ↑ ↓ ☞ ↑ → ↑ ↓ ☞ ☞ → □ □ ☞ □ ☞

→ ← ← ☞ ↑ ↑ → ← ☞ ↑ ☞ □ □ □ □

← ☞ ↑ → → ↑ ↑ ← → → ← → ☞ → □ □

↑ → ↑ → ↑ ☞ ↑ ↓ ↑ ☞ ↑ ↑ ← → → ☞

→ ↑ ↑ → → → → → → → ← ↑ ← → □

↓ ← ↑ ☞ → ↑ → ↑ ☞ ↓ ↓ ↑ ↑ ↓ ☞ □

↑ ☞ ↑ → ↑ ↑ ☞ ☞ ☞ → ← → → □ □ ☞

☞ → ↑ ↑ ↑ ☞ ☞ □ □ ☞ → ↑ ☞ □ □ □

→ ↑ ☞ → ↑ ↓ ↓ ☞ □ ☞ → → ↓ ☞ □ □

↑ ← ☞ ↑ ↓ → → ↑ ↑ ↓ ☞ □ ↓ ↓ ☞

→ → → ↑ ↑ ↓ ↑ ← ☞ ↑ ☞ → → → → ↑

☞ → ← ☞ ↑ ↑ → ↑ ☞ ☞ □ ☞ ↑ ↑ ↑ ↑

← ↑ ← → → ↓ ↓ ☞ ↓ ↓ ☞ → → ↑ ↑ ↑

↑ ↓ ↓ ↓ ☞ ↓ ↓ ☞ ↓ → → ↑ ↑ ☞ → ↑

← → → → → → → → → ↑ ↑ ☞ ↑ ↑

← → → ↑ ↑ ☞ → → → ↑ ☞ ↑ ↑ → ↑ ↑

↑ → → ↑ → → ↑ → ↑ → → ↑ ← ↑ →

← ↑ ↑ ☞ ↑ ↑ ↑ ↑ ↑ ☞ → ↑ → ← ↓ ←

☞ ☞ ↑ → → ↑ ☞ → ↑ → → ↑ ☞ ☞ → →

☞ → ↑ ↑ ☞ ↑ → ↑ ☞ ↑ ↑ ☞ ☞ ☞ ↑ ☞

→ → ↑ ↑ ☞ ↑ ☞ ☞ → ↑ ☞ ↓ ← ↓ ☞ →

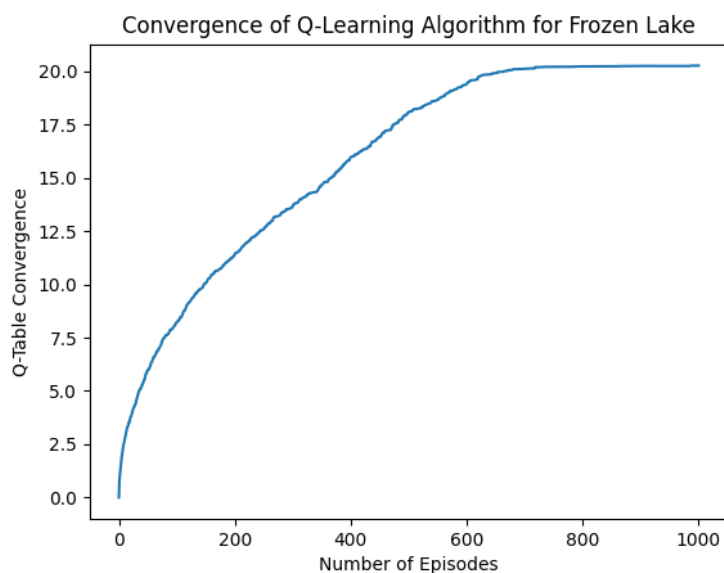
↑ → ↑ → → ↑ → → ↑ ↑ ← ↓ ↓ ← ↑

→ ↑ ↑ ↑ ↑ ☞ → → → ↓ ↓ ↑ ↓ ↓ ☞ ☞

☞ ☞ → → → ↑ ← ☞ ☞ ☞ ↓ ↓ ↑ ↑ ← →

→ → ↑ ↑ ↑ ☞ ↓ → → ↑ ↑ ← → ↑ ↑

نمودار همگرایی که روی مقادیر Q table رسم شده است به شکل زیر می‌باشد. مشاهده می‌شود مقادیر جدول از حدود اپیزود ۸۰۰ به همگرایی می‌رسند.



## 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)

[Frozenlake benchmark - Gymnasium Documentation \(farama.org\)](https://farama.org/environments/frozenlake/)