



گزارش فاز دوم پروژه: تاکسی خودران

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 ϵ -greedy آپدیت می‌کند.) است استفاده شده زیرا نسبت به الگوریتم SARSA که On-policy (مقادیر جدول q_table را مشابه سیاست رفتاری ϵ -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 در دریافت می‌کند. سپس، در یک حلقه به طول حداکثر تعداد دفعاتی که می‌خواهیم عمل یادگیری الگوریتم انجام شود، یک حلقه while تا زمانی که مقدار متغیر done برابر با False باشد می‌زنیم. در این حلقه while با توجه به مشاهداتی که از محیط داریم با استفاده از سیاست ϵ -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 و اطلاعات محیط را برمی‌گرداند.

حال مشاهدات محیط پس از اعمال کنش را به تابع decode داده و موقعیت تاکسی، مسافر و هدف را دریافت می‌کنیم. مقدار مورد انتظار سودمندی در صورت رساندن مسافر به مقصد ۲۰ و در صورت اعمال کنش اشتباه سوار یا پیاده کردن ۱۰- در نظر گرفته شده است. در صورتی که مسافر سوار تاکسی باشد مقدار مورد انتظار سودمندی برابر فاصله اقلیدسی مختصات تاکسی تا مختصات هدف بوده و در صورتی که مسافر سوار تاکسی نباشد مقدار مورد انتظار سودمندی برابر فاصله اقلیدسی مختصات تاکسی تا مختصات مسافر در نظر گرفته می‌شود.

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 برقرار شدن شرط حلقه while از آن 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 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 range(0, MAX_ITER_NUMBER):
4         # reset the environment
5         done = False
6
7         # explore the environment until the agent reaches the goal
8         while not done:
9             # choose the action using epsilon greedy policy
10            action = epsilon_greedy_policy(q_table, observation)
11
12            next_observation, reward, done, info = env.step(action)
13
14            # get the current state of the agent, and the position of the passenger and the goal
15            taxi_row, taxi_col, pas_pos, goal_pos = List(env.decode(next_observation))
16
17            # if the agent reaches the goal, then reward = 20, else if drop and pickup the passenger, then reward = -10
18            if (reward == 20) and done:
19                reward = 20
20            elif (reward == -10) and not done:
21                reward = -10
22            # if the agent picks up the passenger
23            elif (pas_pos == 4) and reward == -1:
24                # Calculate the distance to the goal
25                dist_to_goal = np.linalg.norm(np.array([taxi_row, taxi_col]) - MAP_COLORFUL_PLACES[goal_pos])
26                reward = -0.1 / (1 + np.exp(-dist_to_goal))
27            # if the agent have not pick up the passenger
28            elif (pas_pos != 4) and reward == -1:
29                # Calculate the distance to the passenger
30                dist_to_goal = np.linalg.norm(np.array([taxi_row, taxi_col]) - MAP_COLORFUL_PLACES[pas_pos])
31                reward = -0.1 / (1 + np.exp(-dist_to_goal))
32
33            # choose the next action
34            next_action = np.argmax(q_table[next_observation])
35            # update the q_table
36            q_table[observation][action] = q_table[observation][action] + LEARNING_RATE * \
37            (reward + DISCOUNT_FACTOR * q_table[next_observation][next_action] - q_table[observation][action])
38            # update the observation
39            observation = next_observation
40
41            # if the agent reaches the goal, then reset the environment
42            if done:
43                observation = env.reset()
44
45            # check convergence
46            convergenceTrack.append(np.linalg.norm(q_table.flatten().tolist()))
47            if (i > 1000) and np.isclose(convergenceTrack[-1], convergenceTrack[-2]):
48                print('Values Converged')
49                return
50
51            if (i+1) % 100 == 0:
52                print(i+1)
53
54        print("Training Completed")
55        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 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     map_size = 6
3     with open('q_table_taxi.txt', 'w', encoding="utf-8") as inp:
4         for state in range(map_size**2):
5             x, y = state // map_size, state % map_size
6             if np.array_equal(MAP_COLORFUL_PLACES[0], np.array([x, y])):
7                 inp.write(u'●\t')
8             elif np.array_equal(MAP_COLORFUL_PLACES[1], np.array([x, y])):
9                 inp.write(u'●\t')
10            elif np.array_equal(MAP_COLORFUL_PLACES[2], np.array([x, y])):
11                inp.write(u'●\t')
12            elif np.array_equal(MAP_COLORFUL_PLACES[3], np.array([x, y])):
13                inp.write(u'●\t')
14            else:
15                if np.all(q_table[state] == 0):
16                    inp.write(u'□\t')
17                else:
18                    argm = np.argmax(q_table[state])
19                    if argm == 0:
20                        inp.write(u'↓\t')
21                    elif argm == 1:
22                        inp.write(u'↑\t')
23                    elif argm == 2:
24                        inp.write(u'→\t')
25                    elif argm == 3:
26                        inp.write(u'←\t')
27                    elif argm == 4:
28                        inp.write(u'👤\t')
29                    elif argm == 5:
30                        inp.write(u'🏠\t')
31            if (state + 1) % map_size == 0:
32                inp.write('\n')
```

قسمت main برنامه:

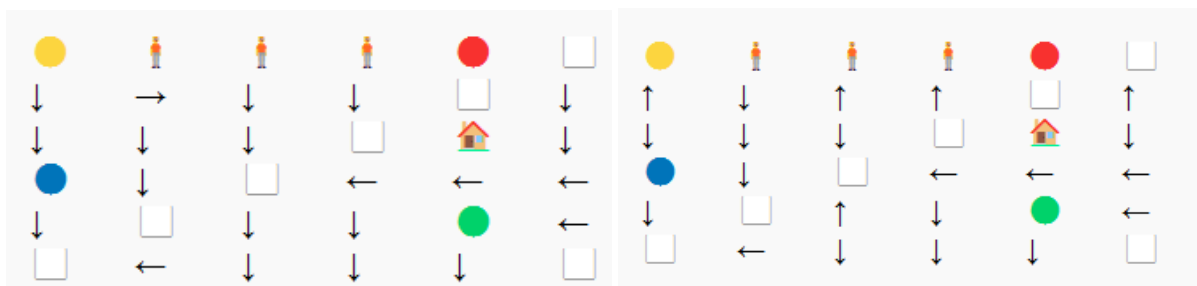
در ابتدا قسمت train را با استفاده الگوریتم Q Learning انجام شده و سیاست به دست آمده را ذخیره می‌کنیم.

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.

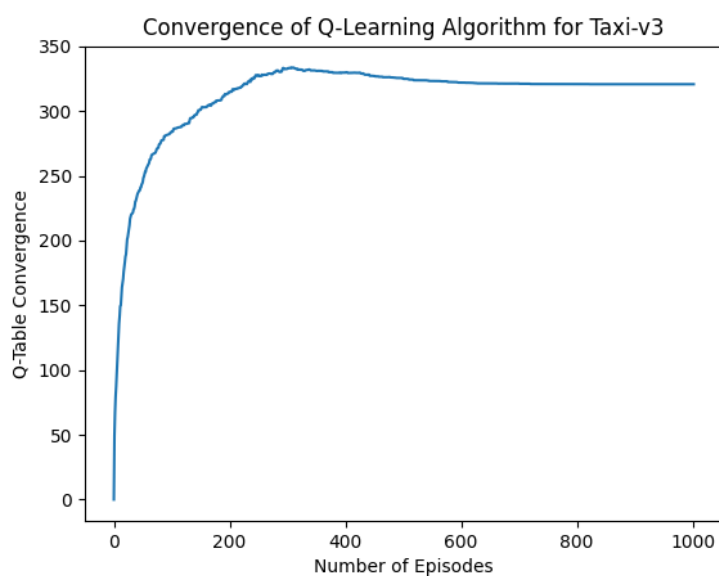
```
1  if __name__ == "__main__":
2
3      # train the agent
4      Q_Learning_Algorithm(observation, q_table)
5
6      save_q_table(q_table)
7
8      # reset the environment
9      observation = env.reset()
10
11     # test the agent for one time
12     done = False
13
14     while not done:
15         action = np.argmax(q_table[observation])
16         observation, reward, done, info = env.step(action)
17
18         print(env.render())
19         sleep(1)
20
21
22     env.close()
23     plot_convergence(convergenceTrack)
```

تحلیل نتایج

با تغییر مقدار اپسیلون از صفر به سمت ۰/۹ مشاهده می‌شود سیاست به دست آمده توسط الگوریتم بهبود پیدا کرده و به سمت سیاست بهینه تغییر پیدا می‌کند. شکل سمت راست مقدار اپسیلون برابر ۰ و شکل سمت چپ مقدار اپسیلون برابر ۰/۹ است. خانه‌های سفید نشان‌دهنده این است که عامل تا به حال در این خانه‌ها نرفته و خانه‌هایی که آدمک دارند نشان‌دهنده انتخاب کنش سوار کردن و خانه ساختمان نشان‌دهنده انتخاب کنش پیاده کردن برای آن حالت از محیط می‌باشد.



نمودار همگرایی که روی مقادیر Q table رسم شده است به شکل زیر می‌باشد. مشاهده می‌شود مقادیر جدول از حدود اپیزود ۶۰۰ 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)

[Solving The Taxi Environment With Q-Learning — A Tutorial | by Wouter van Heeswijk, PhD | Mar, 2023 | Towards Data Science](#)