

گزارش تمرین ۵ام هوش مصنوعی دکتر رهبان

امیر حسین باقری ۹۸۱۰۵۶۲۱

فهرست مطالب

۱	algorithm	۱
۱	code	۲
۳	results	۳
۳	نتیجه تایم train	۱.۳
۳	میانگین و انحراف معیار ۱۰۰۰۰ بار ران تست	۲.۳
۴	نمودار نتایج برای ۱۰۰۰۰ بار تست	۳.۳

algorithm ۱

در اینجا از روش epsilon decay استفاده می کنیم. در این روش مقدار اپسیلون را پس از هر اپیزود کم می کنیم تا explore محیط کمتر شود. برای update مقادیر Q-table مطابق زیر نیز عمل می کنیم.

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t))$$

همچنین برای گسسته سازی فضا بازه -0.07 تا 0.07 که برای سرعت است را به ۱۵ قسمت و بازه 1.2- تا 0.6 که برای بازه مختصات است را به ۱۹ قسمت تقسیم می کنیم.

code ۲

در انتهای هر بخش مانند train و یا تست نتایج را اعم از جدول Q ها و نتیجه ریوارد ها و غیره را در قالب یک npy ذخیره می کنیم.

```
1 class DQAgent:
2     def __init__(self, env):
3         self.env = env
4         num_states = (env.observation_space.high - env.
                        observation_space.low) * np.array([10, 100])
```

```

5         num_states = np.round(num_states, 0).astype(int) + 1
6         self.Q = np.zeros((num_states[0], num_states[1], env.
          action_space.n))
7         self.min = env.observation_space.low
8     def save(self, path, start, end):
9         np.save(path, self.Q)
10        file = open('training_time.txt', 'w')
11        file.write("start time is : "+str(start)+"\n"+"end time is : "+
          str(end)+"\n"+"total time is : "+str(end - start))
12        file.close()
13    def load(self, path):
14        self.Q = np.load(path)
15    def test(self):
16        reward = 0
17        env = self.env
18        Q = self.Q
19        obs = env.reset()
20        done = False
21        while done != True:
22            state = self.get_state(obs)
23            action = np.argmax(Q[state[0], state[1]])
24            obs, r, done, info = env.step(action)
25            reward += r
26            if done and obs[0] >= 0.5:
27                break
28        return reward
29    def get_state(self, obs):
30        return np.round((obs - self.min) * np.array([10, 100])).astype(
          int)
31    def train(self, learning, discount, epsilon, minE, episodes):
32        reduction = (epsilon - minE) / episodes
33        env = self.env
34        Q = self.Q
35        start = datetime.now()
36        for e in range(episodes) :
37            done = False
38            obs = env.reset()
39            current = self.get_state(obs)
40            while done != True:
41                i, j = current[0], current[1]
42                action = np.argmax(Q[i, j]) if np.random.random() >
          epsilon else np.random.randint(0, 3)
43                obs, reward, done, info = env.step(action)
44                next_state = self.get_state(obs)
45                i_p, j_p = next_state[0], next_state[1]
46                if done and obs[0] >= 0.5:
47                    Q[i, j, action] = reward
48                else:
49                    Q[i, j, action] += learning * (reward + discount *
          np.max(Q[i_p, j_p]) - Q[i, j, action])
50                current = next_state
51                epsilon -= reduction
52            env.close()
53            self.Q = Q
54            end = datetime.now()
55            self.save("policy", start, end)
56            return start, end

```

Q-table ها یک جدول $3 \times 19 \times 15$ هستند
بخش train
نیز بر حسب توضیحات داده شده در مقدمه کار می کند.

۳ results

دقت کنید که نتایج حاصل از کد ران شده در دستگاه با نتایج حاصل از کد ران شده در کولب مقدار بسیار اندکی تفاوت بکنند.

۱.۳ نتیجه تایم train

start time is : 2021-07-15 03:02:40.025386
end time is : 2021-07-15 03:10:01.574321
total time is : 0:07:21.548935

۲.۳ میانگین و انحراف معیار ۱۰۰۰۰ بار ران تست

number of more than -140 from 10000 test : 9871 98.71%
number of more than -150 from 10000 test : 10000 100.0%
rewards mean is : -134.9717
rewards std is : 2.8762995515071093

```
main <-  
/home/amirhoosein/Documents/term4/AI/5/RL/venv/bin/python /home/amirhoosein/Documents/term4/AI/5/RL/main.py  
training time : 0:07:21.548935  
number of more than -140 from 10000 test : 9871 98.71%  
number of more than -150 from 10000 test : 10000 100.0%  
rewards mean is : -134.9717  
rewards std is : 2.8762995515071093  
  
Process finished with exit code 0
```

۳.۳ نمودار نتایج برای ۱۰۰۰۰ بار تست

