



دانشگاه زنجان

پروژه:
پیاده‌سازی الگوریتم Hysteretic Q-Learning برای کار
ball-balancing

دانشجو: سینا حیدری

تابستان ۱۴۰۱

مقدمه

در این گزارش، رویکرد Q-Learning hysteric برای یادگیری تقویتی سیستم‌های چندعامله توضیح داده شده است. ما یک آزمایش برای کار ball-balancing پیاده‌سازی کرده‌ایم. در این کار دو عامل که دو سر یک سطح صاف را کنترل می‌کنند باید با action های خود تعادل توپ را بر روی این سطح صاف حفظ کنند. برای روش Q-Learning باید یک فضای حالت گسسته در اختیار داشته باشیم. فضای حالت با توجه به محیط مشخص شده در خود مقاله گسسته شده است. برای مثال، موقعیت توپ بر روی سطح صاف به صد مقدار در بازه ۱- تا ۱ گسسته شده است. در بخش اول این گزارش، آموزش عامل‌ها و در بخش دوم، آزمایش عامل‌های آموزش داده شده بررسی شده است. در بخش سوم نیز روش اجرای پیاده‌سازی و دریافت خروجی‌ها را توضیح داده‌ایم.

۱ آموزش عامل‌ها

فضای حالت و فضای action به شکل زیر گسسته شده است:

- time: Sampling ۰.۳ seconds
- State: Space برای موقعیت صد مقدار بین ۱- و ۱. برای سرعت پنجاه مقدار بین ۳- و ۳.
- action ها پانزده مقدار بین ۱- و ۱.

با این تفاسیر ابعاد Q-table ها به صورت $100 \times 50 \times 15$ خواهد بود. منطق مورد نیاز برای آموزش عامل‌ها در تابع trainHysteretic قرار داده شده است. محتوای این تابع را در پایین مشاهده می‌کنید:

```
1 def trainHysteretic():
2     # create q-Table
3     qTable1 = create_table()
4     qTable2 = create_table()
5     qTables = [qTable1, qTable2]
6
7     iterationRewards = []
8     for i in range(5):
9         rewardSumInTrial = []
10        for trial in range(trials):
11            progress(trial, trials, prefix='Iteration: ' + str(i))
12            # Initialize states (x,xbar)
13            states = (0.495, 1.041)
14            rewardSum = 0
15
16            for t in np.arange(0, 20, samplingTime):
```

```

17     new_actions = nextAction(states, actions, qTables, trial,
18                               numOfEps=40, trials=trials)
19
20     x, v = nextState(h1=new_actions[0], h2=new_actions[1], v=states
21                     [1], t=samplingTime, x_0=states[0],
22                     v_0=states[1])
23
24     # deal with velocities more than 3 and less than -3
25     if v > 3: v = 3
26     if v < -3: v = -3
27
28     if np.abs(x) > 1: break
29
30     thisReward = reward(states[0], states[1])
31     rewardSum = rewardSum + thisReward
32
33     new_states = (np.round(x, decimals=decimals), np.round(v,
34                     decimals=decimals))
35     # check if the new states are in discretized states
36     new_states = isValidState(new_states)
37
38     qTables = hysteretic(qTables, states, new_actions, alpha, beta,
39                           thisReward, gamma, new_states)
40     states = new_states
41
42     rewardSumInTrial.append(rewardSum)
43
44     iterationRewards.append(rewardSumInTrial)
45     mean_output = np.mean(iterationRewards, axis=0)
46     plt.plot(list(range(trials)), mean_output, '-')
47     plt.title('Hysteretic')
48     plt.savefig('./Plots/Hysteretic.png')
49     plt.clf()
50
51     pickle.dump(qTables[0], open('QTables/qT1_Hysteretic.p', 'wb'))
52     pickle.dump(qTables[1], open('QTables/qT2_Hysteretic.p', 'wb'))

```

در این تابع در ابتدا دو نمونه Q-table (یکی برای هر عامل) ساخته شده است. در حلقه‌ای که در خط ۱۰ مشاهده می‌شود، منطق مربوط به هر trial جا داده شده است. در هر trial که بیست ثانیه طول می‌کشد (sampling time=0.03) عامل‌ها یک action انجام می‌دهند و آن action توپ را به حالت بعدی می‌برد. انتخاب action با استفاده از روش ϵ -greedy انجام شده است. دینامیک مورد نیاز برای به دست آوردن حالت بعدی با استفاده از معادله زیر به دست می‌آید:

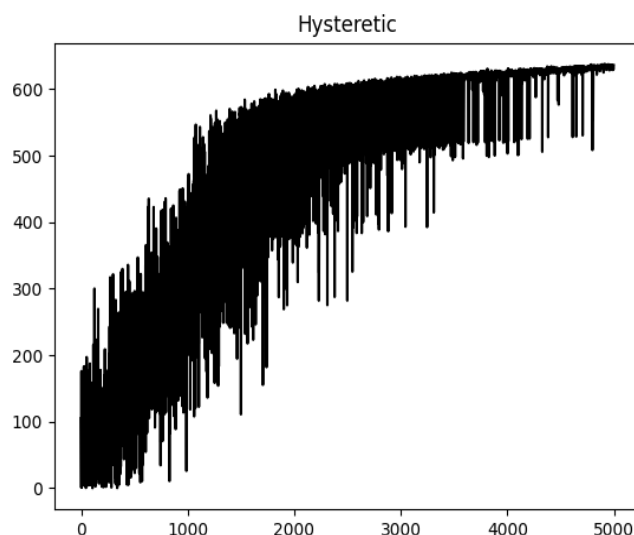
$$m\ddot{x} = -c\dot{x} + mg\left(\frac{h_1 - h_2}{l}\right)$$

سپس با توجه حالت جدید (موقعیت و سرعت) عامل‌ها reward دریافت می‌کنند. reward برای هر حالت با توجه به فرمول زیر به دست می‌آید:

سپس با توجه به reward دریافت شده و الگوریتم Q-Learning hysteretic جدول Q-table برای

$$r = 0.8e^{-\frac{x^2}{0.25^2}} + 0.2e^{-\frac{\dot{x}^2}{0.25^2}}$$

هر عامل به روزرسانی می‌شود. reward در هر واحد زمانی محاسبه شده و مجموع آن برای هر trial ذخیره می‌شود. تعداد trial ها مثل مقاله ۵۰۰۰ در نظر گرفته شده است. ما پنج بار این تعداد trial ها را تکرار می‌کنیم و از آن میانگیری می‌کنیم (دقیقا مثل مقاله با این تفاوت که در مقاله بیست بار تکرار شده است). نمودار تشویق‌ها بعد از آزمایش را در تصویر زیر مشاهده می‌کنید:



همانطور که در تصویر بالا مشاهده می‌کنید در trial های بعدی reward بیشینه شده است. این یعنی الگوریتم Q-Learning hysteretic کار اصلی‌اش که maximization reward است را درست انجام داده است.

۲ آزمایش عامل‌ها

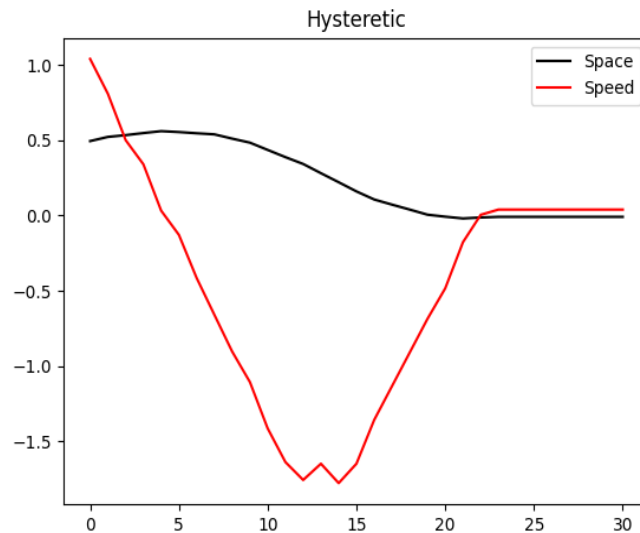
منطق تست عامل‌ها در تابع test قرار داده شده است. کد این تابع را در پایین مشاهده می‌کنید:

```

1 def test(qTables):
2     spaces, velocities = [], []
3     states = (0.495, 1.041)
4     spaces.append(states[0])
5     velocities.append(states[1])
6
7     for i in range(30):
8         actions = nextAction(states, None, qTables, None)
9         x, v = nextState(h1=actions[0], h2=actions[1], v=states[1], t=0.03,
10                        x_0=states[0], v_0=states[1])
11         spaces.append(x)
12         velocities.append(v)
13         # if the ball falls from the flat surface
14         if np.abs(x) > 1:
15             break
16         new_states = (np.round(x, decimals=3), np.round(v, decimals=3))
17         new_states = isValidState(new_states)
18         states = new_states
19
20     plt.plot(spaces, '-', label="Space", color='black')
21     plt.plot(velocities, '-', label="Speed", color='red')
22     plt.legend()
23     plt.title("Hysteretic")
24     plt.savefig('./Plots/' + "Hysteretic" + '_test.png')
25     plt.clf()

```

برای آزمایش عامل‌ها همانند مقاله ابتدا از یک وضعیت اولیه شروع کرده و در ۳۰ قدم عامل با توجه حالت یک تصمیم می‌گیرد. با توجه به تصمیم عامل وضعیت توپ تغییر می‌کند. نمودار مربوط به آزمایش عامل‌ها را در تصویر زیر مشاهده می‌کنید:



همانطور که مشاهده می‌کنید، به مرور زمان توپ در وسط سطح صاف و با سرعت پایین balance شده

است.

۳ اجرای پیاده‌سازی

برای اجرای پیاده‌سازی ابتدا در پوشه code دستور زیر را اجرا کنید تا پکیج‌های مورد نیاز نصب شوند:

```
pip install -r requirements.txt
```

سپس برای آموزش عامل‌ها دستور زیر وارد شود:

```
python training.py
```

و برای آزمایش عامل‌ها دستور زیر وارد شود:

```
python test.py
```