

بسمه تعالی



دانشگاه تحصیلات تکمیلی علوم پایه زنجان

گروه علوم رایانه و فناوری اطلاعات

درس طراحی سیستم های هوشمند



تمرین چهارم

Maze

محمد عمید عباسی 924468 - محمد سروش فرخزاده 924430 - مسعود فتحی 924429

بهار 96

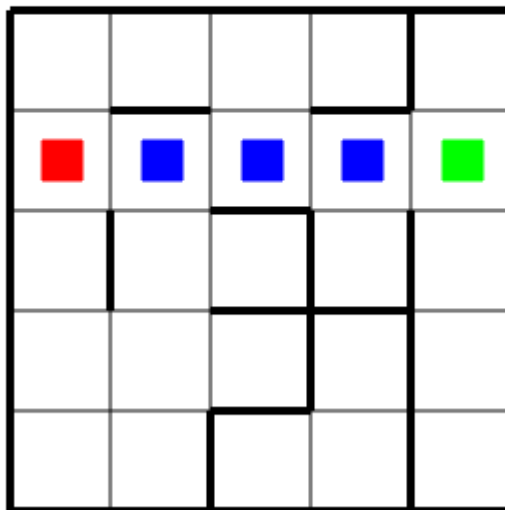
فهرست

3	مقدمه.....
4	الگوریتم value iteration.....
6	الگوریتم policy iteration.....
9	الگوریتم یادگیری Q.....
11	ابتدای هر m file
12	انتهای هر m file
13	مقایسه ی الگوریتم های value و policy
14	مشکلات الگوریتم Q.....
16	پاسخ سوالات.....

مقدمه

این پروژه با استفاده از تکنولوژی matlab پیاده سازی شده و شامل شش m file می باشد که برای اجرای الگوریتم value iteration فایل maze.m ، برای اجرای الگوریتم policy iteration فایل maze_policy.m و برای اجرای الگوریتم Q فایل Q2.m اجرا گردد؛ سه m file بعدی توابع استفاده شده در برنامه می باشند که در ادامه توضیح داده می شوند.

هدف از اجرای هر سه الگوریتم value iteration ، policy iteration و Q به دست آوردن تابع سیاست بهینه می باشد که بیانگر بهترین حرکت در هر حالت می باشد.



الگوریتم value iteration

Value iteration یکی از الگوریتم های حوزه ی یادگیری تقویتی در یادگیری ماشین می باشد که در آن هدف به دست آوردن تابع سیاست می باشد. در این الگوریتم توابع reward و sigma برای عامل شناخته شده می باشد.

1. For each state s , initialize $V(s) := 0$.
2. Repeat until convergence {
For every state, update $V(s) := R(s) + \max_{a \in A} \gamma \sum_{s'} P_{sa}(s') V(s')$.
}

This algorithm can be thought of as repeatedly trying to update the estimated value function using Bellman Equations (2).

Value iteration algorithm

در این الگوریتم ابتدا value های همه ی خانه ها با صفر مقداردهی اولیه می شوند، سپس تا زمانی که همگرایی رخ دهد، value هر خانه طبق فرمول آپدیت می شود.

منظور از همگرایی در این الگوریتم یکسان بودن ماتریس value در دو iteration متوالی می باشد.

فرمول value به این صورت می باشد که در هر دور برای هر خانه مقدار مجموع reward و γ ضربدر value خانه ی بعد را به ازای همه ی action های ممکن آن خانه حساب می کند و مقدار ماکزیمم را جایگزین value فعلی خانه می کند.

پیاده سازی این الگوریتم به صورت زیر می باشد:

```

29 while ~(all(G2 == G1))
30     G1 = value;
31     for i = 1:(s*s)
32         if any(goal == i)
33             value(i)=0;
34         else
35             if B(i,1)
36                 temp_val(1) = reward(i,'N',goal,s) + gama*value(sigma(i,'N',s));
37             end
38             if B(i,2)
39                 temp_val(2) = reward(i,'E',goal,s) + gama*value(sigma(i,'E',s));
40             end
41             if B(i,3)
42                 temp_val(3) = reward(i,'S',goal,s) + gama*value(sigma(i,'S',s));
43             end
44             if B(i,4)
45                 temp_val(4) = reward(i,'W',goal,s) + gama*value(sigma(i,'W',s));
46             end
47             [value(i),policy(i)] = max(temp_val);
48             temp_val = zeros(4,1);
49         end
50     end
51     G2 = value;
52     iter = iter + 1;
53
54 end

```

الگوریتم policy iteration

policy iteration یکی از الگوریتم های حوزه ی یادگیری تقویتی در یادگیری ماشین می باشد که در آن هدف به دست آوردن تابع سیاست می باشد. در این الگوریتم توابع reward و sigma برای عامل شناخته شده می باشد.

Policy Iteration

1. Set $n = 0$ and initialize the policy π^0 so that it has non negative expected reward.

2. Let v^n be the solution to the system of equations

$$v(i) = r(i) + \sum_j p_{ij}^{\pi^n} v(j)$$

where $p_{ij}^{\pi^n}$ is the probability of moving from state i to state j under the policy π^n . When multiple solutions are found, the minimal non-negative one is preferred.

3. For each state s with action space $A(s)$, set

$$\pi^{n+1}(s) \in \operatorname{argmax}_{a \in A(s)} \sum_j p_{ij}^a v^n(j)$$

breaking ties so that $\pi^{n+1}(s) = \pi^n(s)$ whenever possible.

4. If $\pi^{n+1}(s) = \pi^n(s)$ for all s , stop and set $\pi^* = \pi^n$. Otherwise increment n by 1 and return to step 2.

□ π^0 is initialized to the near optimal policy s^*

□ It is an alternative way to find an optimal policy π^*

□ States in poker tournaments are made of vectors containing the stack size for each player

■ e.g. 3-players state = (x1, x2, x3)

■ x1: stack of the button

■ x2: stack of the small blind

■ x3: stack of the big blind

Policy iteration algorithm

در این الگوریتم ابتدا مقدار policy برای همه ی خانه ها به صورت رندوم مقداری اولیه می شود و سپس تا زمانی که همگرایی رخ دهد دو عمل زیر تکرار می شود:

عمل یک: به دست آوردن value هر خانه با توجه به policy مرحله ی قبل. که به دو صورت قابل پیاده سازی می باشد:

1-روش ریاضی: در این روش با استفاده از policy مرحله ی قبل معادله ای از value ها به دست می آوریم و معادله را حل می کنیم. برای حل این معادله می توان از linear regression استفاده کرد. که فرمول آن به صورت زیر می باشد:

$$X = (A^T A)^{-1} * A^T * B$$

که در این فرمول ماتریس X ماتریس value ها و ماتریس A ، ماتریس ضرایب و B ماتریس اعداد طرف دیگر معادله می باشد.

2- روش استفاده از simplified value iteration: در این روش با استفاده از value iteration ساده شده مقدار value ها را با توجه به policy های مرحله ی قبل به دست می آوریم؛ که ما در این برنامه از این روش استفاده کرده ایم.

تفاوت Value iteration ساده شده با value iteration اصلی در این می باشد که در value iteration اصلی در هر iteration ما کمترین را انتخاب می کنیم ولی در value iteration ساده شده در هر iteration ، value را با توجه به policy مرحله ی قبل انتخاب می کنیم و نه مقدار ماکزیمم را!!!!

عمل دوم: به روز رسانی مقدار policy می باشد. که در این عمل policy یا action انتخاب می شود که مقدار زیر را max کند:

$$\text{Reward}(s,a) + \text{gama} * v(\text{sigma}(s,a))$$

پیاده سازی:

```
31 while ~(all(G2 == G1))
32     G1 = policy;
33     for i = 1:(s*s)
34         for j = 1:(s*s)
35             if any(goal == j)
36                 value(j) = 0;
37             else
38                 if B(j,policy(j))
39                     value(j) = reward(j,policy_ac(policy(j)),goal,s) + gama*value(sigma(j,policy_ac(policy(j)),s));
40                 end
41             end
42         end
43         if any(goal == i)
44             value(i)=0;
45         else
46             if B(i,1)
47                 temp_val(1) = reward(i,'N',goal,s) + gama*value(sigma(i,'N',s));
48             end
49             if B(i,2)
50                 temp_val(2) = reward(i,'E',goal,s) + gama*value(sigma(i,'E',s));
51             end
52             if B(i,3)
53                 temp_val(3) = reward(i,'S',goal,s) + gama*value(sigma(i,'S',s));
54             end
55             if B(i,4)
56                 temp_val(4) = reward(i,'W',goal,s) + gama*value(sigma(i,'W',s));
57             end
58             [~,policy(i)] = max(temp_val);
59             temp_val = zeros(4,1);
60         end
61     end
62     G2 = policy;
63     iter = iter + 1;
64 end
```


الگوریتم یادگیری Q

Q یکی از الگوریتم های حوزه ی یادگیری تقویتی در یادگیری ماشین می باشد که در آن هدف به دست آوردن تابع سیاست می باشد. در این الگوریتم توابع reward و sigma برای عامل شناخته شده نمی باشد و عامل فقط sigma و reward لحظه ای را می داند.

Q learning algorithm

For each s, a initialize the table entry $\hat{Q}(s, a)$ to zero.

Observe the current state s

Do forever:

- Select an action a and execute it
- Receive immediate reward r
- Observe the new state s'
- Update the table entry for $\hat{Q}(s, a)$ as follows:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

- $s \leftarrow s'$
-

Q-learning algorithm

در این الگوریتم یک جدول Q به اندازه تعداد عمل ها ضربدر تعداد حالت ها می سازیم و آن را با صفر مقدار دهی اولیه می کنیم سپس یک حالت رندوم را انتخاب می کنیم و وارد حلقه ی بی نهایت می شویم:

یک عمل را انتخاب می کنیم و reward ناشی از آن را از محیط دریافت می کنیم و حالت s' ناشی از آن را مشاهده می کنیم؛ سپس جدول Q را مطابق فرمول به روز رسانی می کنیم و وارد state جدید می شویم. زمانی که به حالت goal رسیدیم یک episode تمام می شود. و برای episode جدید دوباره الگوریتم را تکرار می کنیم.

پیاده سازی:

```
28
29 for i = 0:maxiter
30     %ka = ka + 1;
31     % if i>250000
32     %     ka = 10;
33     % end
34     if(any(current_state == goal))
35         if(h~=1)
36             for u = 1:h-1
37                 Q_table(w(1,u)) = w(2,u);
38             end
39             h = 1;
40         end
41         current_state = randi([1 s*s],1);
42     else
43         for l = 1 : 4
44             aa = aa + ka*Q_table((current_state-1)*4+1);
45         end
46         for l = 1:4
47             kaa(l) = ka*Q_table((current_state-1)*4+1)/aa;
48         end
49         action =(2*eps)*floor(rand/(2*eps));
50         aa = 0;
51
52         if (0 <= action) && (action < kaa(1))
53             action = 1;
54         elseif (kaa(1) <= action) && (action < (kaa(1) + kaa(2)))
55             action = 2;
56         elseif (((kaa(1) + kaa(2)) <=action) && (action < (kaa(1) + kaa(2) + kaa(3))))
57             action = 3;
58         else
59             action = 4;
60         end
61         %action = randi([1 4],1);
62         if(B(current_state,action))
63             w(1,h) = ((current_state - 1)*4) + action;
64             w(2,h) = reward(current_state,policy_ac(action),goal,s) + ...
65             gama*max(Q_table(((sigma(current_state,policy_ac(action),s)-1)*4 + 1):((sigma(current_state,policy_ac(action),s)-1)*4)+ 4));
66             current_state = sigma(current_state,policy_ac(action),s);
67             h = h + 1;
68         elseif (g>2000)
69             current_state = randi([1 s*s],1);
70             g = 0;
71             h = 1;
72         else
73             g = g + 1;
74         end
75     end
76 end
```

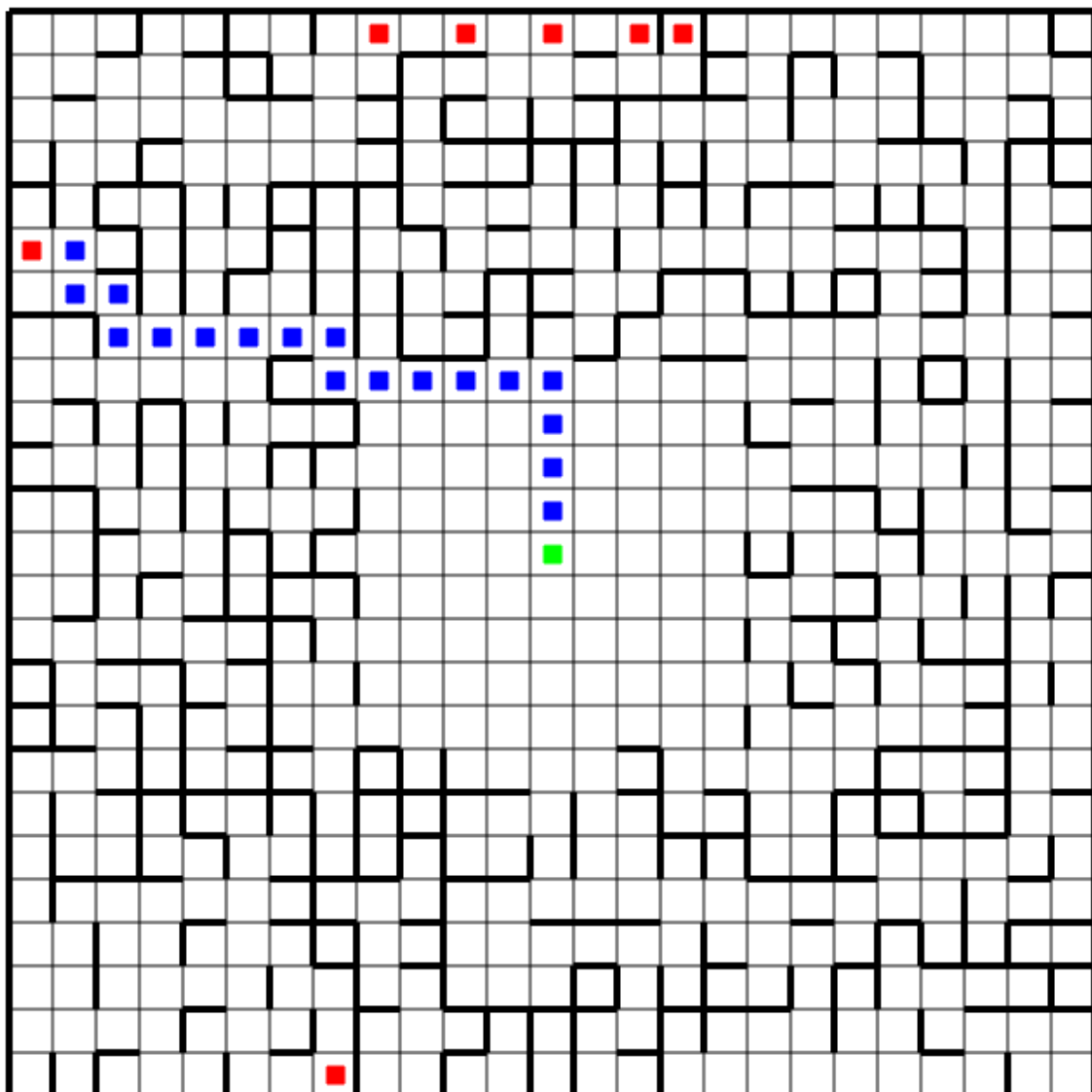
ابتدای هر m file

در ابتدای هر m file، فایل های ورودی برنامه خوانده می شوند و در یک آرایه ذخیره می شوند.

```
1 |fileid = fopen('maze_25.txt','r');
2 A = fscanf(fileid,'%u');
3 s = A(1); %size of matrix
4 if s == 5
5     start = A(103);
6     goal = A(104:end);
7     A = A(3:102);
8     B = reshape(A,[4,25]);
9     B = B';
10 elseif s == 25
11     start = A(2503);
12     goal = A(2504:end);
13     A = A(3:2502);
14     B = reshape(A,[4,625]);
15     B = B';
16 end
```

انتهای هر m file

در انتهای هر m file، با استفاده از policy به دست آمده maze مربوطه کشیده می شود؛ به این صورت که با توجه به فایل ورودی خط های عمودی و افقی ترسیم می گردد. سپس با توجه به نقطه ی شروع و پایان(ها) و policy مسیر شروع تا پایان کشیده می شود.



مقایسه ی الگوریتم های value و policy

بدون در نظر گرفتن بحث محاسباتی هر دو الگوریتم اثبات همگرایی دارند و هر دو به جواب بهینه همگرا می شوند.

الگوریتم value iteration در هر iteration فقط همسایه های خانه هایی را که مقدار دارند به روز می کند، در صورتی که policy iteration در هر iteration تمام خانه ها را به روز رسانی می کند، که این باعث می شود در بعضی از مسائل در iteration های کمتری به جواب بهینه برسد.

Policy تا حدودی به مقدار دهی اولیه هم وابستگی دارد.

در کل برای مسائل کوچک policy ممکن است بهتر از value عمل کند ولی برای مسائل بزرگ value بهتر است چون policy ممکن است به جواب نرسد.

مشکلات الگوریتم Q

1-بزرگ بودن جدول Q .

2- رندوم انتخاب کردن عملها.

راه حل: استفاده از راه حل احتمالاتی. با استفاده از این روش احتمال انتخاب خانه هایی که مقدار آن ها کمتر می باشد یعنی آپدیت نشده اند را در اوایل اجرا که نیاز به explore داریم و احتمال انتخاب خانه هایی را که مقدار بیشتری دارند یعنی آپدیت شده اند را در اواخر اجرا که نیاز به explit داریم افزایش می دهیم.

K عددی است که در اوایل مسئله مقدارش کم و رفته رفته زیاد می شود.

فرمول استفاده شده :

Exploration versus Exploitation

- The Q-learning algorithm doesn't say how we could choose an action
- If we choose an action that maximises our estimate of Q we could end up not exploring better alternatives
- To converge on the true Q values we must favour higher estimated Q values but still have a chance of choosing worse estimated Q values for exploration (see the convergence proof of the Q-learning algorithm in [Mitchell, sec. 13.3.4.]). An action selection function of the following form may employed, where $k>0$:

$$P(a_k | s) = \frac{k^{\hat{Q}(s, a_k)}}{\sum_{k'} k^{\hat{Q}(s, a_{k'})}}$$

3-در هر iteration فقط یک خانه رو آپدیت می کند.

راه حل: همه ی خانه های یک episode را ذخیره و زمانی که به goal رسید update می کنیم بدین صورت تعداد iterationها را کاهش می دهیم

پاسخ سوالات

الف) بله چون در هر قدم عملی که value را بیشینه می کند، انتخاب می گردد. Value هم برابر است با جمع ریوارد های تخفیف داده شده تا هدف.

ب) مانند قسمت الف.

ج) در قسمت مقایسه ی دو الگوریتم policy و value ذکر گردید.

چ) مانند الف.

د) در قسمت مشکلات الگوریتم Q ذکر گردید.