



به نام خدا



دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر
سیستم های هوشمند

تمرین سری 5

شایان واصف	نام و نام خانوادگی
810197603	شماره دانشجویی
23 دی ماه	تاریخ ارسال گزارش

فهرست گزارش سوالات

- سوال 1 – عنوان سوال 3
- سوال ۲ – عنوان سوال 5
- پیاده سازی : 5
-: $9.0 = \gamma$ 6
-: $1 = \gamma$ 7
- تاثیر افزایش هزینه : 8
- تحلیلی : 10
- سوال 3 – عنوان سوال 12
- امتیازی : 16

سوال 1 – عنوان سوال

در ابتدا دیتاست داده شده را load می‌کنیم و سپس دقت را روی مجموعه Train و Test به ترتیب اعلام می‌کنیم :

Train data

```
▶ b = []  
for i in mush_train.index:  
    #print(classify(mush.loc[i,features]),mush.loc[i,target])  
    b.append(classify(mush_train.loc[i,features]) == mush_train.loc[i,target])  
print(sum(b),"correct of",len(mush_train))  
print("Accuracy:", sum(b)/len(mush_train))
```

📄 7563 correct of 7600
Accuracy: 0.9951315789473684

شکل 1-1 : دقت در دادگان آموزش

Test data

```
▶ #Test data  
b = []  
for i in mush_test .index:  
    #print(classify(mush.loc[i,features]),mush.loc[i,target])  
    b.append(classify(mush_test .loc[i,features]) == mush_test .loc[i,target])  
print(sum(b),"correct of",len(mush_test ))  
print("Accuracy:",sum(b)/len(mush_test ))
```

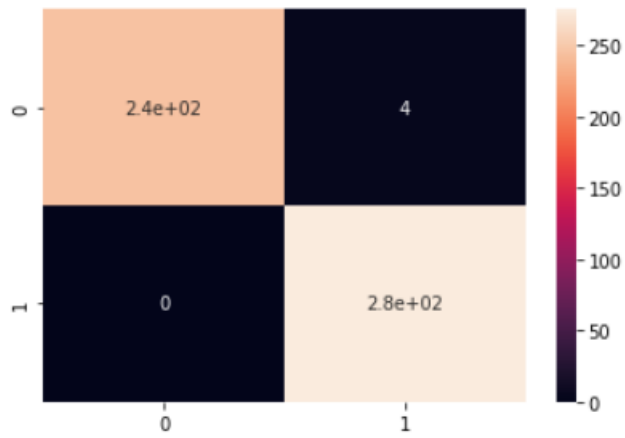
📄 520 correct of 524
Accuracy: 0.9923664122137404

شکل 2-1 : دقت در دادگان تست

در نهایت طبق خواسته سوال ، ماتریس آشفته را رسم می کنیم :

```
sns.heatmap(A,annot=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3134442d90>
```



شکل 1-3: ماتریس آشفته در داده‌گان تست

سوال ۲ – عنوان سوال

پیاده سازی :

در ابتدا اطلاعات داده شده توسط مساله را در زیر خلاصه می کنیم :

Initialization

Maximum

- MAX_Capacity = 20
- MAX_transfer = 5

FIRST of the appointment

- A_REQUEST_FIRST = 3
- B_REQUEST_FIRST = 4

END of the appointment

- A_REQUEST_END = 3
- B_REQUEST_END = 2

gamma factor(default 0.9)

- DISCOUNT = 0.9

Define reward

- CREDIT = 10
- COST = 2

شکل 1-1-2 : خلاصه اطلاعات حاصل از مدل سازی

در ادامه از سه اصطلاح زیر در کد بهره می بریم :

Return states

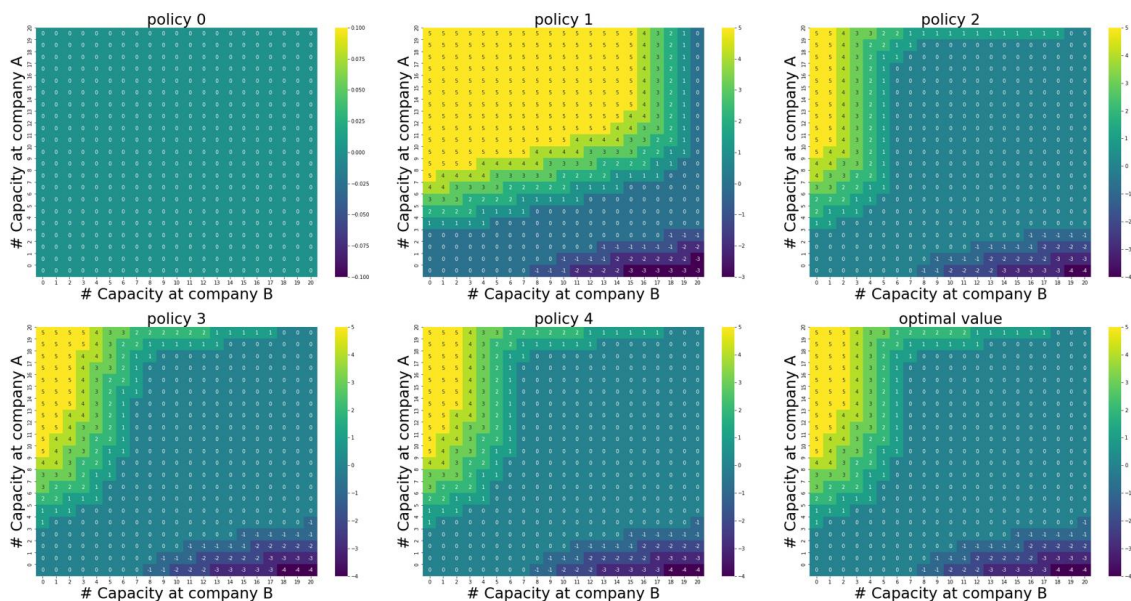
Overall info

- state: [# of capacity in first company , # of capacity in second company]
- action:
 - positive if transferin from first company to second one,
 - negative if transferring from second company to frist one,
- stateValue: state value matrix

شکل 2-1-2 : تعریف اطلاعات استفاده شده در کد

$$\gamma = 0.9$$

نمودار سیاست های اتخاذ شده بر حسب تعداد ظرفیت شرکت A و B به کمک کتابخانه seaborn و دستور sns.heatmap در زیر رسم شده است :

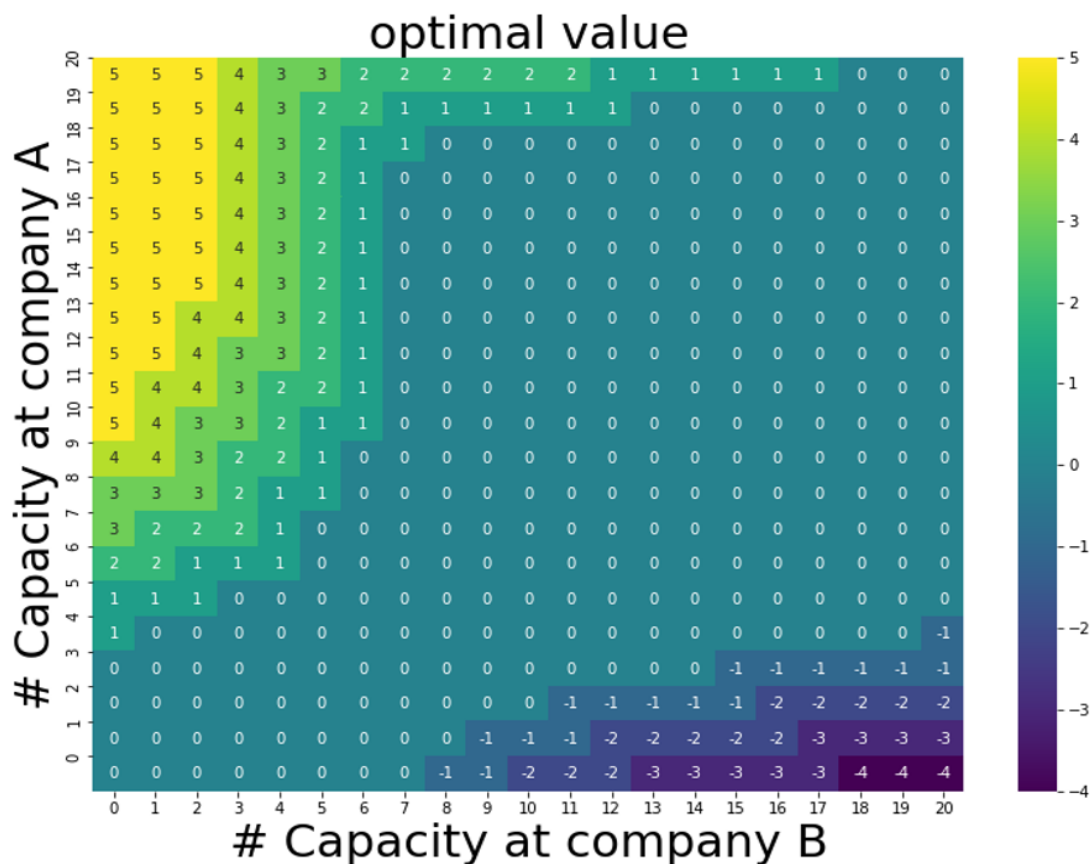


شکل 2-1-3 : heat map حاصل از سیاست بدست آمده در هر تکرار ($\gamma = 0.9$)

طبق نمودار های بالا ، بعد از 6 بار بروز رسانی تابع V به همراه سیاست اتخاذ شده ، در 5 امین تکرار به سیاست بهینه می‌رسیم . اعداد نشان داده شده در هر heat map ، سیاست متناظر با ظرفیت شرکت A و B را نشان می‌دهند.

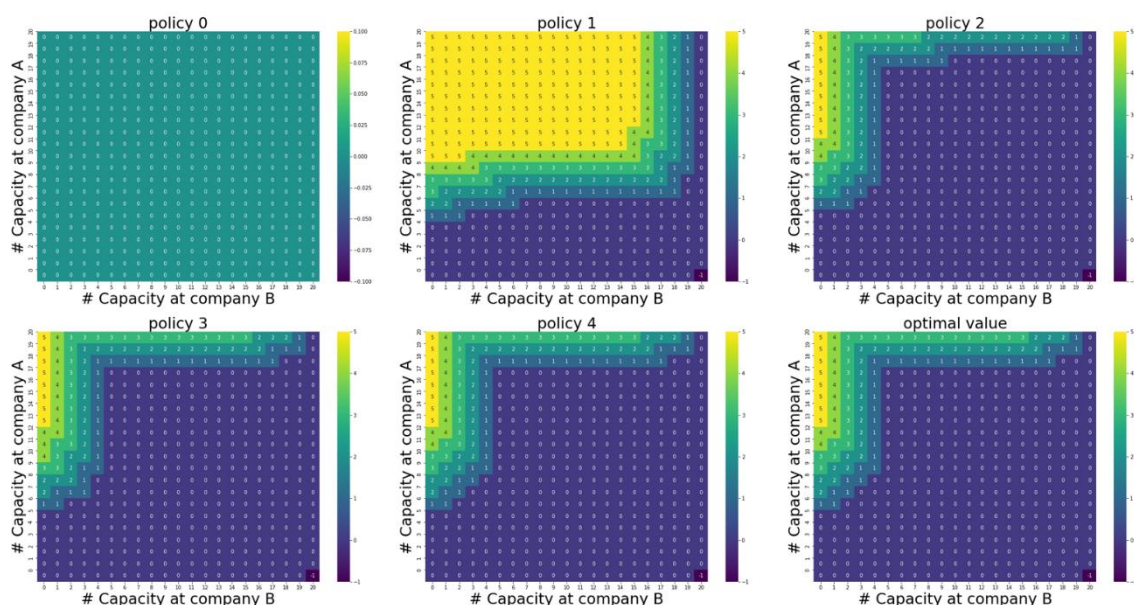
برای مثال ، در صورتی که شرکت B ، تنها 2 ظرفیت داشته باشد و شرکت A بزرگتر مساوی 14 عدد ظرفیت داشته باشد ، سیاست بهینه این خواهد بود که شرکت B از تمام ظرفیت خود برای خرید هر 5 ظرفیت شرکت A استفاده کند .

نمای دقیق تر سیاست بهینه نهایی به صورت زیر می باشد :



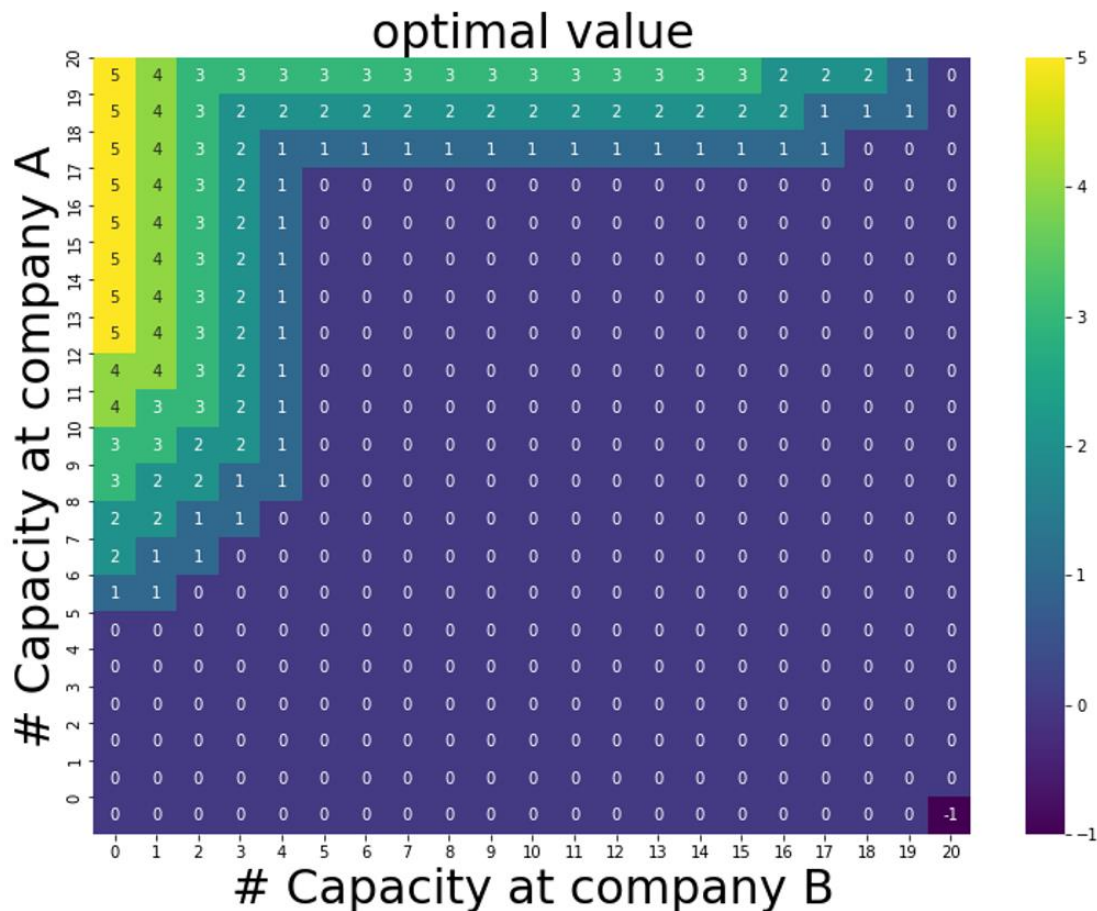
شکل 2-1-4 : heat map حاصل از سیاست بهینه همگرا شده ($\gamma = 0.9$)

$\gamma = 1$ ☹



شکل 2-1-5 : heat map حاصل از سیاست بدست آمده در هر تکرار ($\gamma = 1$)

نمای دقیق تر سیاست بهینه نهایی به صورت زیر می باشد :

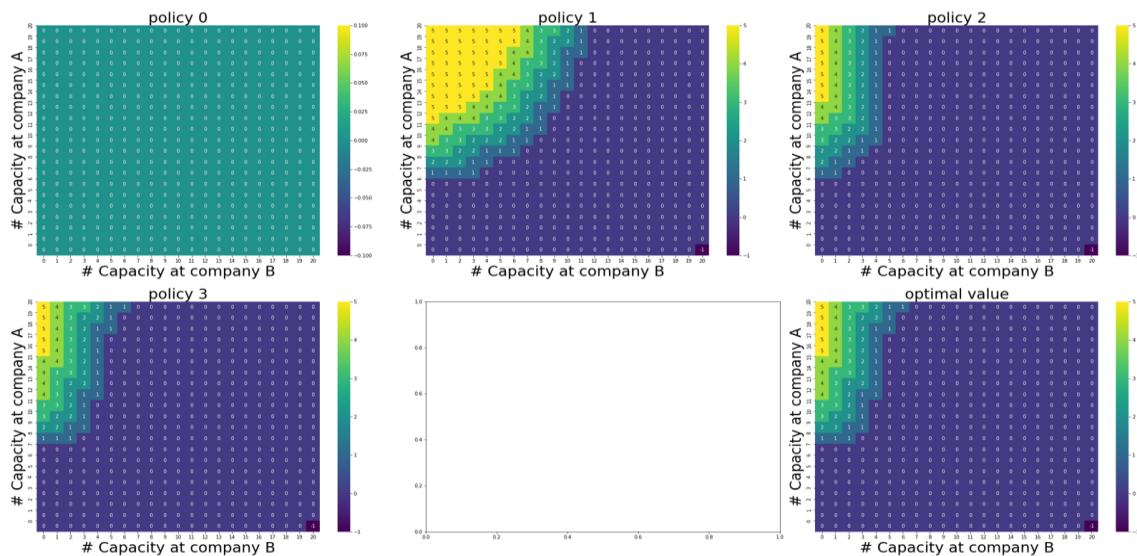


شکل 2-1-4 : heat map حاصل از سیاست بهینه همگرا شده ($\gamma = 0.9$)

طبق نمودار های بدست آمده برای ضریب تخفیف برابر 1 ، تعداد سیاست های مبنی بر انتقال بین دو شرکت کاهش می باید . همچنین شرکت A تنها زمانی که تعداد ظرفیت آن برابر صفر باشد ، 1 واحد از شرکت B خریداری می کند .

تاثیر افزایش هزینه :

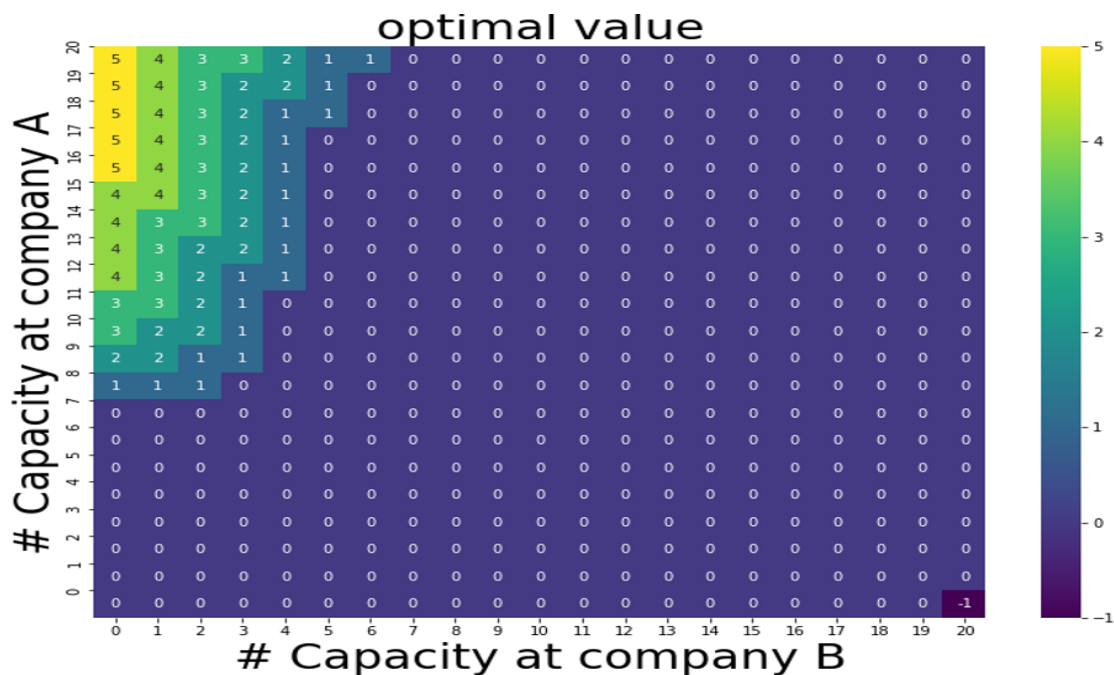
طبیعتاً افزایش هزینه با $\gamma = 1$ ، باید مانند افزایش γ هنگامی که هزینه ها ثابت است ، عمل کند. چرا که در رابطه bellman ، هر دو باعث افزایش تابع هزینه مربوط به حالات جدید می شوند. ادعای گفته شده را در زیر نیز با رسم نمودار ها می توان بررسی کرد :



شکل 5-1-2: heat map حاصل از سیاست بدست آمده در هر تکرار ($Cost = 6$)

طبق نمودار بالا ، با افزایش هزینه از 2 به 4 ، همگرایی از حالت $\gamma = 1$ هم سریع تر شده ولی سیاست های اتخاذ شده در هر دو تقریباً یکسان می باشند . (با کمی دقت می توان متوجه شد که تعداد سیاست های برابر صفر ($action=0$) در این حالت بیشتر نیز می باشد) .
در واقع گویی با افزایش ضرر و ایجاد پاداش منفی بزرگتر ، دو شرکت ترغیب کمتری نسبت به خرید یا فروش کالا به هم دارند چون در هر صورت یکی از آنها ضرر زیادی متحمل خواهد شد.

نمای دقیق تر سیاست بهینه نهایی به صورت زیر می باشد :



شکل 6-1-2: heat map حاصل از سیاست بهینه همگرا شده ($Cost = 6$)

تحلیلی:

شماره واهف - ۸۱۰۱۹۷۶.۳ بنام صرا تیرن ۵- سوال ۲- تحلیلی

طبق مسئله انجام شد... سرفه یادگیری بود. در این مسئله - کار می آید و $\gamma = 1$ می باشد.

طبق الگوریتم Policy iteration، ابتدا یک سیاست رندم در نظر می گیریم و سپس تابع ارزش را آپدیت می کنیم و این کار مدام انجام می دهیم،

در هر مرحله به معادله ی بلوچن، تابع ارزش را برزسانی می کنیم.

سیاست اولیه

←	↑	←	3
→	↓	→	-2
↑	■	↑	↑

$$V_{\pi}(s) = R_s^{\pi(s)} + \gamma \sum_{s'} P_{ss'}^{\pi(s)} V_{\pi}(s'), \quad R_s^{\pi(s)} = 0, \quad \gamma = 1$$

در ابتدا در خانه ای که در مجاور خانه های با یادداشت 3 و 2 - هستند را آپدیت می کنیم.

$$V(3,3) = 0 + 1 \left(\frac{1}{6} V(2,3) + \frac{1}{2} V(3,2) + \frac{1}{2} V(4,3) \right) = 0.6$$

$$V(3,2) = 0 + 1 \left(\frac{1}{6} V(4,2) + \frac{1}{2} V(3,3) + \frac{1}{2} V(3,1) \right) = -1.2$$

V-function

0	0	0.6	3
0	0	-1.2	-2
0	■	0	0

$$\begin{cases} V(2,3) = 0 + 1 \left(\frac{1}{6} V(2,3) + \frac{1}{2} V(3,3) + \frac{1}{2} V(1,3) \right) \\ \rightarrow \frac{1}{4} V(2,3) = \frac{1}{4} \rightarrow V(2,3) = 0.3 \\ V(1,1) = 0 + 1 \left(\frac{1}{6} V(4,2) + \frac{1}{2} V(4,1) + \frac{1}{2} V(3,1) \right) \\ \rightarrow \frac{1}{8} V(4,1) = -1.2 \rightarrow V(4,1) = -1.5 \end{cases}$$

0	0.3	0.6	3
0	0	-1.2	-2
0	■	0	-1.5

$$\begin{cases} V(2,2) = 0 + 1 \left(\frac{1}{6} V(2,2) + \frac{1}{2} V(1,2) + \frac{1}{2} V(3,2) \right) \\ \rightarrow \frac{1}{4} V(2,2) = -1.24 \rightarrow V(2,2) = -1.6 \\ V(3,1) = 0 + 1 \left(\frac{1}{6} V(3,2) + \frac{1}{2} V(4,1) + \frac{1}{2} V(3,1) \right) \\ \rightarrow \frac{1}{8} V(3,1) = -1.72 - 1.3 = -1.02 \rightarrow V(3,1) = -1.125 \end{cases}$$

0	0.3	0.6	3
0	-1.6	-1.2	-2
0	■	-1.125	-1.5

$$\begin{cases} V(1,3) = 0 + 1 \left(\frac{1}{6} V(1,3) + \frac{1}{2} V(1,3) + \frac{1}{2} V(1,2) \right) \rightarrow V(1,3) = 0 \\ V(1,2) = 0 + 1 \left(\frac{1}{6} V(2,2) + \frac{1}{2} V(1,3) + \frac{1}{2} V(1,1) \right) = -1.36 \\ V(1,1) = 0 + 1 \left(\frac{1}{6} V(1,2) + \frac{1}{2} V(1,1) + \frac{1}{2} V(1,1) \right) \rightarrow \frac{1}{6} V(1,1) = -1.276 \\ \rightarrow V(1,1) = -1.36 \end{cases}$$

0	-1.36	0.6	3
-1.36	-1.6	-1.2	-2
-1.36	■	-1.125	-1.5

policy improvement

→	→	→	3
↑	↑	↑	-2
↑	■	↑	←

انجام Iteration اول

شکل 1-2-2: محاسبه Value function و Policy در اولین تکرار

برای تکرار دوم خواهم داشت

$$V(3,3) = 0 + 1 \left(\underset{.136}{.16 \times 3} + \underset{-.14}{.12 V(3,3)} + \underset{-.12}{.12 \times -112} \right) \rightarrow .18 V(3,3) = .16 + 118 - .174 = 21.22$$

$$V(3,2) = 0 + 1 \left(\underset{.136}{.16 \times .16} + \underset{-.14}{.12 \times -2} + \underset{-.12}{.12 \times -.16} \right) \rightarrow V(3,2) = -.176$$

$$V(2,3) = 0 + 1 \left(\underset{.136}{.16 \times .16} + .12 \times V(2,3) + \underset{-.12}{.12 \times -.16} \right) \rightarrow .18 V(2,3) = .154 \rightarrow V(2,3) = .375$$

$$V(4,1) = 0 + 1 \left(.16 \times -1128 + .12 \times -2 + .12 \times V(4,1) \right) \rightarrow .18 V(4,1) = -21668 \rightarrow V(4,1) = -1835$$

$$V(2,2) = 0 + 1 \left(.16 \times .13 + .12 \times -.136 + .12 \times -112 \right) \rightarrow V(2,2) = -.1732$$

$$V(3,1) = 0 + 1 \left(.16 \times -112 + .12 \times V(3,1) + .12 \times -1835 \right) \rightarrow .18 V(3,1) = -11.87 \rightarrow V(3,1) = -1.358$$

$$V(1,3) = 0 + 1 \left(.16 \times .375 + .12 \times V(1,3) + .12 \times -.136 \right) \rightarrow .18 V(1,3) = .153 \rightarrow V(1,3) = .17$$

$$V(1,2) = 0 + 1 \left(.16 \times .17 + .12 \times V(1,2) + .12 \times -.1732 \right) \Rightarrow .18 V(1,2) = .1876 \rightarrow V(1,2) = .777$$

$$V(1,1) = 0 + 1 \left(.16 \times .17 + .12 \times V(1,1) + .12 \times V(1,1) \right) \rightarrow .16 V(1,1) = .166 \rightarrow V(1,1) = .777$$

.17	.138	21.22	3
.777	-.173	-.176	-2
.777		-1.36	-1835

policy
Improvement

→	→	→	3
↑	↑	↑	-2
↑		↑	←

optimal policy after 2 iterations

شکل 2-2-2: محاسبه Value function و Policy در دومین تکرار

سوال 3 – عنوان سوال

در ابتدا طبق کد داده شده ، محیط بازی که در آن پاداش ها و موانع تعریف شده اند را پیاده سازی می کنیم . در ادامه به کمک صدا زدن کلاس GridWorldEnv ، تعداد حالت ها (States) و اعمال (actions) و در نتیجه ماتریس Q را به صورت تمام یک تشکیل می دهیم :

```
env = GridworldEnv()
action_space_size = env.action_space.n
state_space_size = env.observation_space.n
Q_table = np.ones((state_space_size, action_space_size))
```

```
print(action_space_size)
print(state_space_size)
print(Q_table.shape)
```

```
4
100
(100, 4)
```

شکل 4-1 : ایجاد محیط بازی و تشکیل ماتریس Q

دلیل این موضوع که از ماتریس تمام 1 به جای تمام صفر به عنوان مقدار دهی Q استفاده کردیم ، تسریع در فرآیند همگرایی الگوریتم می باشد . طبیعتاً ، در مجموع 100 حالت و 4 عمل مختلف را می توانیم انجام دهیم .

در ادامه Hyper-parameter ها را تعریف می کنیم :

```
num_of_episodes = 20000
max_steps_per_episode = 200

learning_rate = 0.7
discount_rate = 0.99

exploration_rate = 1
```

شکل 4-2 : مقدار دهی Hyper-parameter های محیط بازی

در زیر هر کدام از پارامتر های بالا را توضیح می دهیم :

☛ Number of episodes: هر episode ، به یک دور حرکت Agent در محیط شبیه سازی شده می گویند. ما در هر episode ، Agent رو در جای اول قرار داده و باتوجه به ماتریس Q آپدیت

شده در مرحله قبل ، reward های جدید را بدست آورده و مشابه قبل ماتریس Q را آپدیت می کنیم .

☛ Maximum steps per episode : Agent ما در هر episode تنها قادر است تعداد محدودی گام بردارد و بعد از آن دوباره به جای اول باز میگردد.

☛ Learning rate : LR نرخ بروز رسانی ماتریس Q می باشد که تعیین می کند که به چه مقدار از q مرحله قبل استفاده شود و چه مقدار از مقادیر بدست آمده تاثیر پذیرد و برای مثال اگر LR برابر 1 باشد ، ماتریس Q جدید کاملاً مستقل از ماتریس Q در مرحله قبل بدست می آید .

☛ Discount rate : همان γ می باشد که بر روی مقدار ماکسیمم q مرحله جدید تاثیر می گذارد.

☛ Exploration rate : در هر مرحله از حرکت Agent ، مساله ای بنام Exploration Vs Exploitation مطرح می باشد.

Exploitation از دانسته های ما نشأت می گیرد و در واقع ما بر دانسته های خود تکیه می کنیم. Exploration از نادانسته ها می آید و مربوط به کشف حالت های ناشناخته و ریسک پذیری می باشد. بنابراین اگر در تمامی مراحل تنها بر دانسته های خود تکیه کنیم ، امکان دارد به جواب بهینه اصلی همگرا نشویم و برای مثال امکان تعداد پاداش های بازی زیاد باشد و به یک local optimum همگرا شویم .

در ابتدا که Agent چیزی از اطلاعات اطراف نمی داند ، Exploration را برابر 1 قرار می دهیم و به مرور زمان ، با افزایش دانسته های خود از محیط ، Exploration کاهش می یابد و Exploitation افزایش می یابد . به میزان کاهش Exploration ، $\text{Exploration_decay_rate}$ می گوئیم که یا می توان آن را دستی تنظیم کرد و یا می توان آنرا معکوس episode تعریف کرد که ما از این روش استفاده می کنیم . در واقع با افزایش episode ، دانش ما از محیط اطراف افزایش می باشد و بنابراین نیاز ما به Exploration کمتر می باشد . پس با افزایش Episode ، Exploration باید کاهش یابد که به صورت زیر انجام می شود :

```
exploration_rate = 1 - np.log(episode + 1) / np.log(num_of_episodes + 1)
```

شکل 3-4 : نحوه تعریف $\text{Exploration_decay_rate}$

بنابراین در هر مرحله یک عدد رندوم بین صفر و یک از توزیع uniform انتخاب می کنیم و آنرا با Exploration آن مرحله مقایسه می کنیم . در صورتی که مقدار آن عدد از Exploration بیشتر باشد ، Action ای را انتخاب می کنیم که متناظر با بیشترین مقدار q باشد و در عوض اگر کمتر باشد ، یک action تصادفی را انتخاب می کنیم که کد مربوط به آن در زیر آمده است :

```

random_number = random.uniform(0,1)

if random_number > exploration_rate:
    action = np.argmax(Q_table[state,:])

else:
    action = env.action_space.sample()

```

شکل 4-4 : نحوه انتخاب action به کمک Exploration-rate

بعد از توضیحات بالا ، به هسته اصلی کد یعنی آپدیت ماتریس Q می پردازیم که با پارامتر های بالا به صورت زیر تعریف می شود :

$$q^{new}(s, a) = (1 - \alpha) \underbrace{q(s, a)}_{\text{old value}} + \alpha \overbrace{\left(R_{t+1} + \gamma \max_{a'} q(s', a') \right)}^{\text{learned value}}$$

شکل 4-4 : نحوه بروز رسانی ماتریس Q به کمک پارامتر های داده شده

که در کد به صورت زیر ایجاد می کنیم :

```

Q_table[state, action] = Q_table[state, action] * (1 - learning_rate) \
+ learning_rate * (reward + discount_rate * np.max(Q_table[new_state,:]))

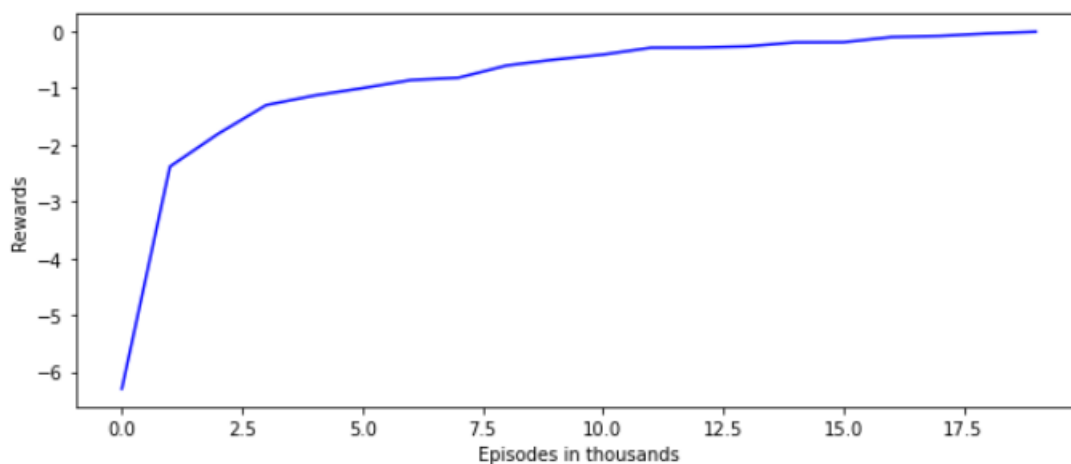
```

شکل 4-4 : نحوه بروز رسانی ماتریس Q در کد

در نهایت پاداش را بر حسب هر 1000 تعداد Episode خروجی می گیریم و رسم می کنیم :

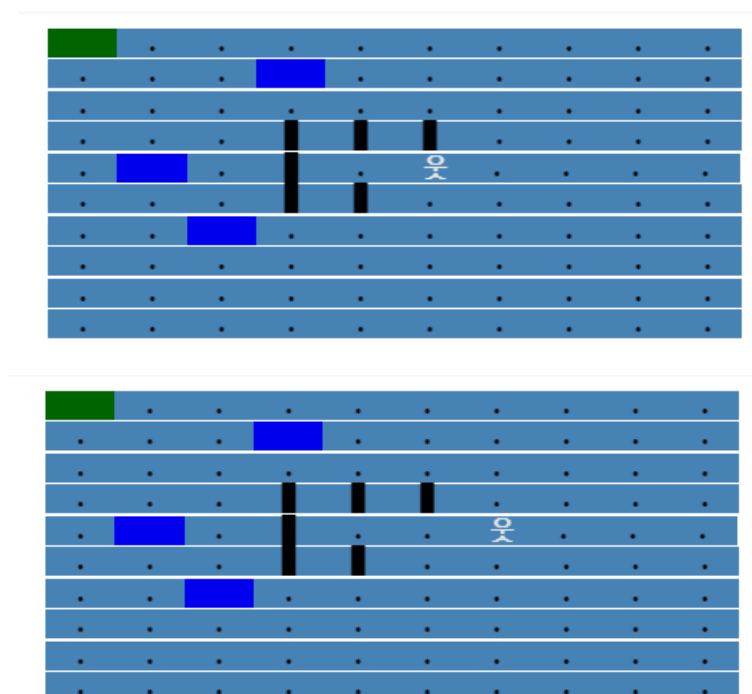
Average reward per thousand episodes	11000 : -0.408
	12000 : -0.288
1000 : -6.295	13000 : -0.285
2000 : -2.381	14000 : -0.263
3000 : -1.805	15000 : -0.193
4000 : -1.298	16000 : -0.19
5000 : -1.13	17000 : -0.1
6000 : -1.0	18000 : -0.081
7000 : -0.858	19000 : -0.037
8000 : -0.816	20000 : -0.006
9000 : -0.601	
10000 : -0.496	

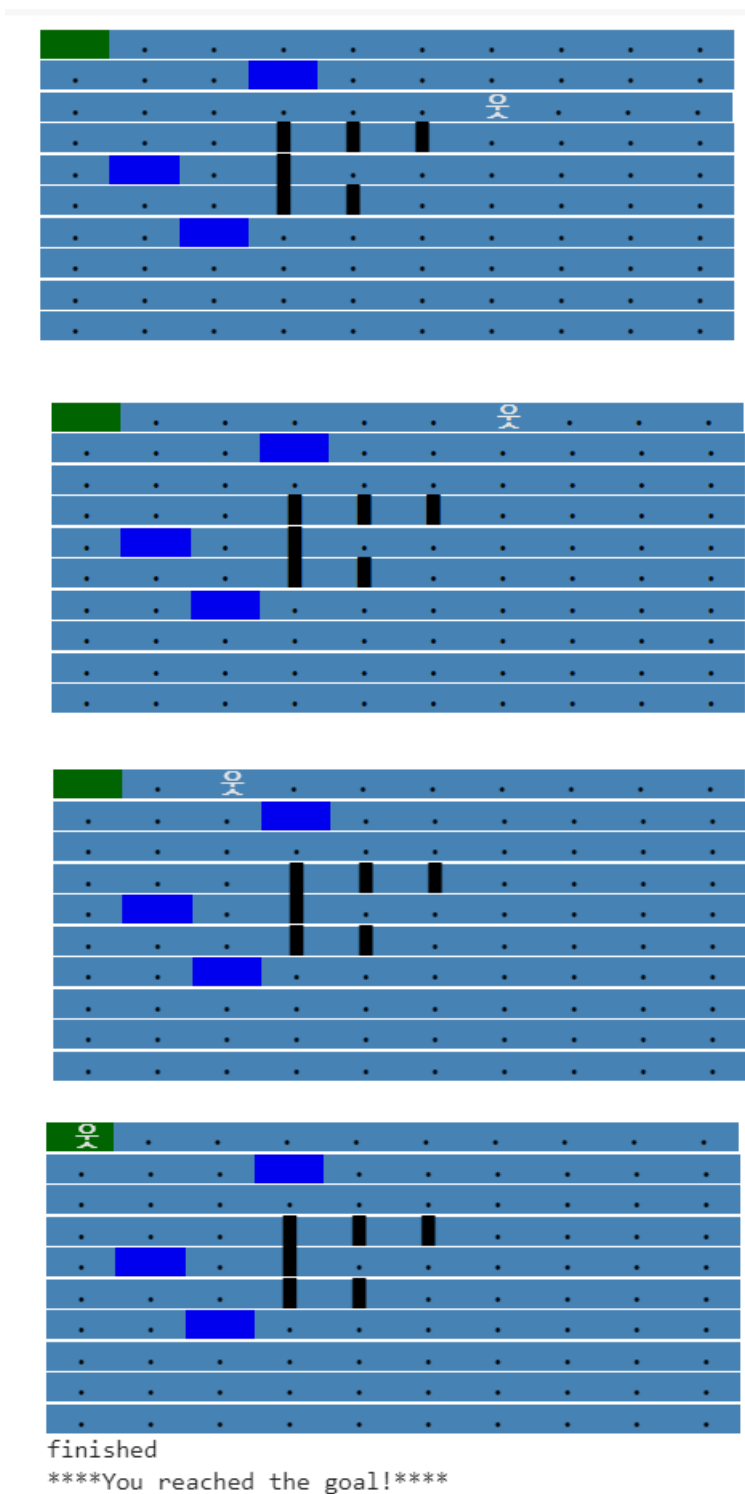
شکل 4-5 : میانگین پاداش بدست آمده در هر 1000 episode



شکل 4-6 : نمودار میانگین پاداش بدست آمده در هر 1000 episode

طبق خروجی بالا ، میانگین پاداش بدست آمده از یک عدد منفی بزرگ به سمت صفر نزدیک می شود تا در نهایت در صفر پایدار می شود . بنابراین طبق پاداش های داده شده ، در بهترین حالت ، agent ما پاداش منفی دریافت نمی کند و در واقع یاد می گیرد که به صخره ها و مکان های پر عمق نرود و همچنین در نهایت به جزیره برسد . در زیر روند حرکت agent و حالت نهایی توقف او را نشان می دهد (لازم به ذکر است که agent در حرکت های خود طبق آخرین ماتریس آپدیت شده Q تصمیم می گیرد) :





شکل 4-7: نحوه حرکت agent تا رسیدن به جزیره

امتیازی:

طبق قسمت قبل ، agent به جزیره همگرا شد ولی قصد داریم تا سرعت آنرا کمی بهبود بخشیم . بنابراین نیاز است تا انگیزه ای برای agent ایجاد کنیم تا به سمت جزیره برود . اگر به ابتدای نمودار ها

نگاه کنید ، پیک مهمی در ابتدای کار می خورد که به نظر می رسد اگر این پیک بیشتر باشد ، agent در تعداد Episode کمتری محیط را یاد می گیرد. بنابراین چون در ابتدای کار با صخره ها برخورد داریم و بعد به مناطق عمیق می رسیم ، پاداش منفی مربوط به این دو را با هم جابجا می کنیم . از طرفی در تعدادی نقاط ، پاداش نسبتا کوچکی قرار می دهیم تا میانگین پاداش ها در هر Episode همچنان سه سمت صفر نزدیک شود (به جزیره همگرا شویم) ولی تا حدی agent را در رسیدن به جزیره بیشتر کمک کند . در کد زیر نحوه ایجاد پاداش های مثبت و منفی را آورده ایم :

```
#Rewards of each state
reward = 1.0 if is_done(s) else 0.0
reward = -103.0 if pit(s) else reward
#Generating rewards
pre_help1=lambda s :s in [40,41,42,28,48]
pre_help2=lambda s :s in [10]
reward = -1 if pre_help1(s) else reward
reward = 0.01 if pre_help2(s) else reward
```

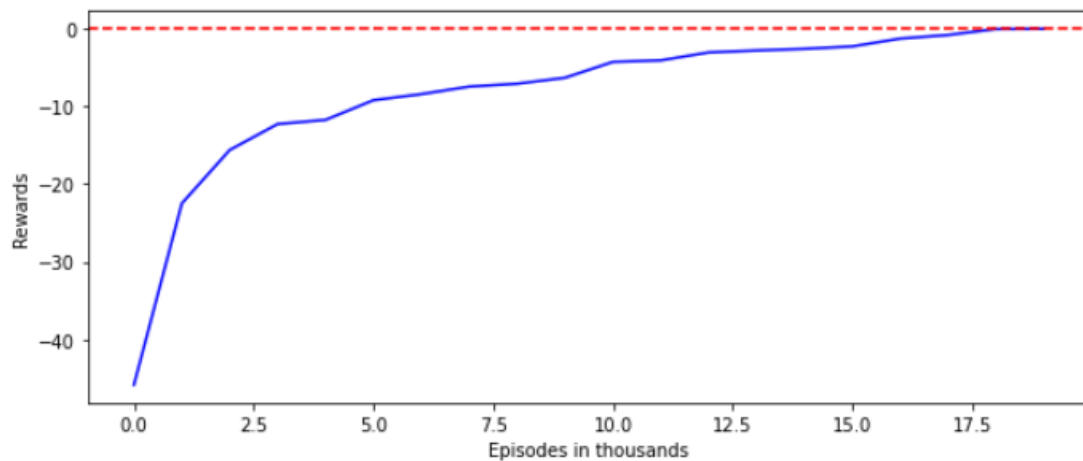
شکل 4-8 : نحوه ایجاد پاداش منفی و مثبت در حالت جدید

با اجرای الگوریتم ، میانگیم مقادیر پاداش بر حسب هر 1000 Episode به صورت زیر می باشد :

	Average reward per thousand episodes
11000 : -4.14129999999999	1000 : -52.276979999999774
12000 : -4.4002799999999903	2000 : -21.466010000000013
13000 : -3.65836999999998918	3000 : -15.3637799999999847
14000 : -2.66419999999999646	4000 : -15.5906699999999834
15000 : -2.50421999999999735	5000 : -11.866829999999984
16000 : -2.6431499999999923	6000 : -8.4017999999999852
17000 : -1.6391600000000008	7000 : -7.2844299999999881
18000 : -1.48311000000000086	8000 : -6.7984799999999858
19000 : -0.7180000000000007	9000 : -7.5254199999999853
20000 : -0.163010000000000318	10000 : -5.8104299999999886

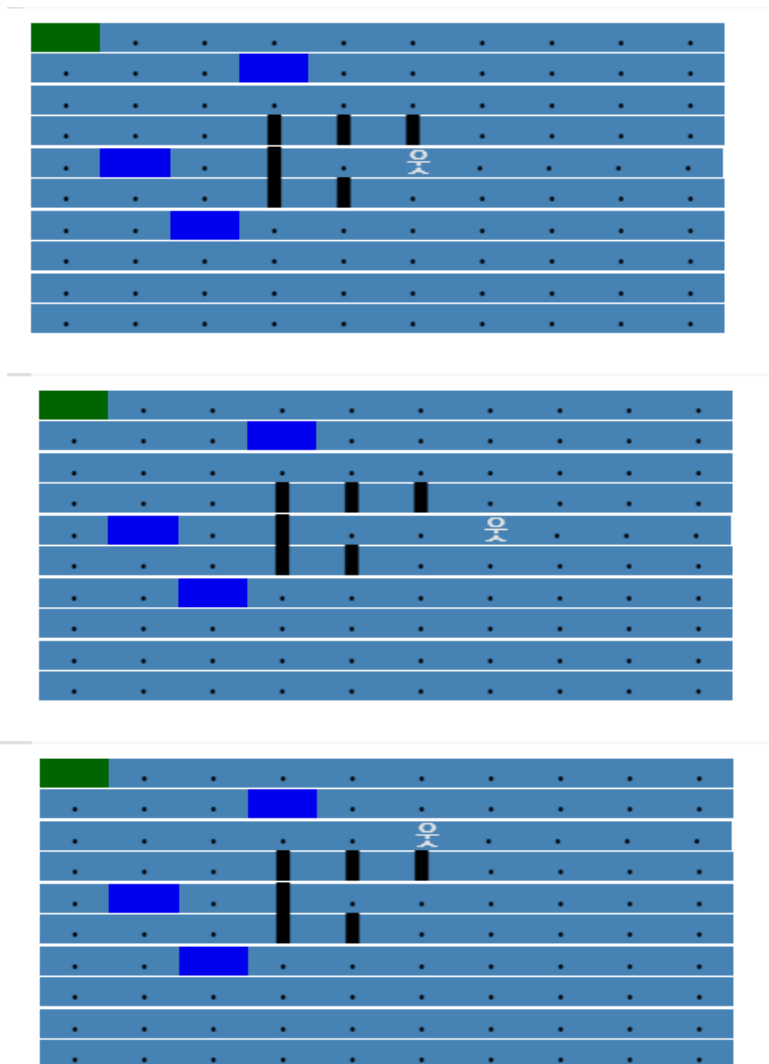
شکل 4-9 : مقادیر میانگین پاداش بر حسب Episode

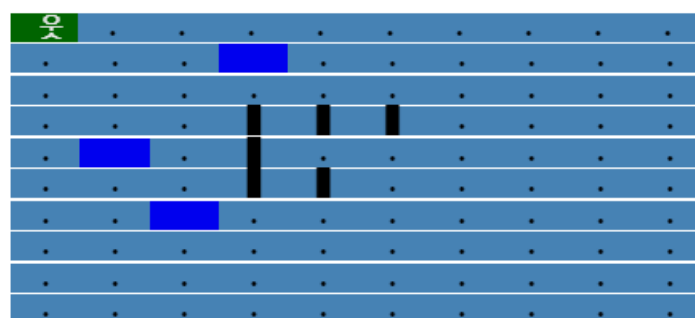
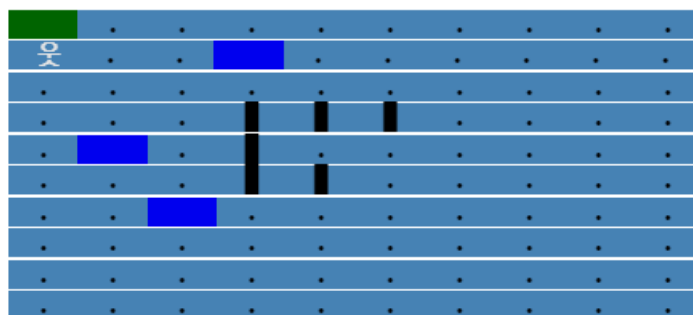
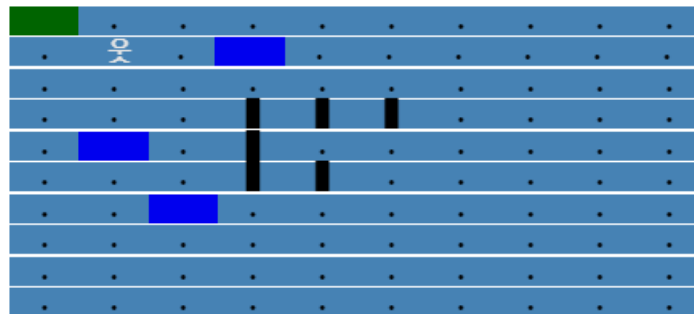
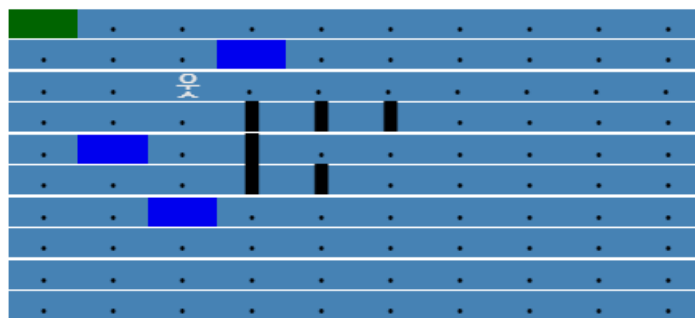
و نمودار پاداش بر حسب Episode به صورت زیر می باشد :



شکل 4-10 : نمودار پاداش بر حسب Episode با پاداش های جدید

حال در زیر ، نحوه جابجایی agent را در محیط رسم می کنیم (لازم به ذکر است که agent در حرکت های خود طبق آخرین ماتریس آپدیت شده Q تصمیم می گیرد) :





finished
 ****You reached the goal!****

شکل 8-4: نحوه حرکت agent تا رسیدن به جزیره

همانطور که مشاهده می‌شود، در این حالت با تغییر در روند پاداش ها، agent نسبت به قبل از اطلاعات بیشتری نسبت به محیط برخوردار است و در نتیجه با سرعت بیشتری همگرا می‌شود.