



یادگیری تقویت شده

گزارش پروژه پنجم

نام استاد : دکتر ستوده

نام دانشجو : حمزه قاندي

شماره دانشجویی : 9831419

تاریخ : 99/5/10

| | |
|---------|--|
| 2..... | بخش اول |
| 3..... | مقدمه |
| 4..... | یادگیری تقویت شده |
| 5..... | Q-function: |
| 5..... | Q-Learning: |
| 5..... | Greedy action |
| 6..... | ε-Greedy action |
| 6..... | پیاده سازی انتخاب با احتمال دلخواه |
| 6..... | Deep Q-Learning (DQN): |
| 6..... | overestimation |
| 7..... | Double DQN(DDQN) |
| 7..... | بخش دوم |
| 7..... | تشریح مسئله |
| 8..... | شبکه های عصبی به عنوان تقریبی از تابع سیاستگذاری [Policy function] |
| 11..... | نقش مکانیزم فراموشی : |
| 12..... | نتیجه گیری |

تاکنون، با دو پارادایم کلی در مسائل یادگیری ماشین^۱ آشنا شده ایم، یادگیری بانظارت^۲ و یادگیری بدون نظارت^۳. در این پروژه، با سومین پارادایم در یادگیری ماشین، به نام یادگیری تقویت شده^۴ آشنا میشویم. در زیر مختصراً به تشریح این سه دسته میپردازیم

SUPERVISED LEARNING



UNSUPERVISED LEARNING



REINFORCEMENT LEARNING



یادگیری با نظارت :

متداول ترین پارادایم در یادگیری ماشین میباشد. در این حالت ما تعدادی ورودی^۵ به همراه خروجی^۶ (یا برچسب) نظیر مربوط به هر ورودی را در اختیار داریم و معنی یادگیری در این پارادایم، یافتن نوعی رابطه یا نگاشت بین ورودیها و خروجیهای نظیرشان است. مسائل این حوزه نیز خود به دو دسته رگرسیون^۷ و طبقه بندی^۸ تقسیم شده اند

یادگیری بدون نظارت:

در این روش، ورودی ها را در اختیار داریم اما خروجی نظیرشان را نداریم (داده ها برچسب ندارند) و معنی و هدف یادگیری در این حالت، یافتن الگوها، ساختارها و همبستگی بین این داده هاست (مسائل خوشه بندی^۹) و یا در حالت کلیتر یافتن توزیع احتمالی که داده ها از آن استخراج شده اند.

یادگیری تقویت شده:

در این پارادایم، معمولاً داده ای (چه به صورت برچسب دار و یا بدون برچسب) نداریم و عامل یادگیر طی تعامل با محیط اطراف خود، اطلاعاتی را جمع آوری کرده و بر اساس این اطلاعات اقداماتی را در جهت نیل به هدف (معمولاً کسب بیشترین جایزه) انجام میدهد در ادامه به تشریح این پارادایم یادگیری میپردازیم.

¹ Machine Learning

² Supervised

³ Unsupervised

⁴ Reinforcement Learning

⁵ Input

⁶ Output/Label

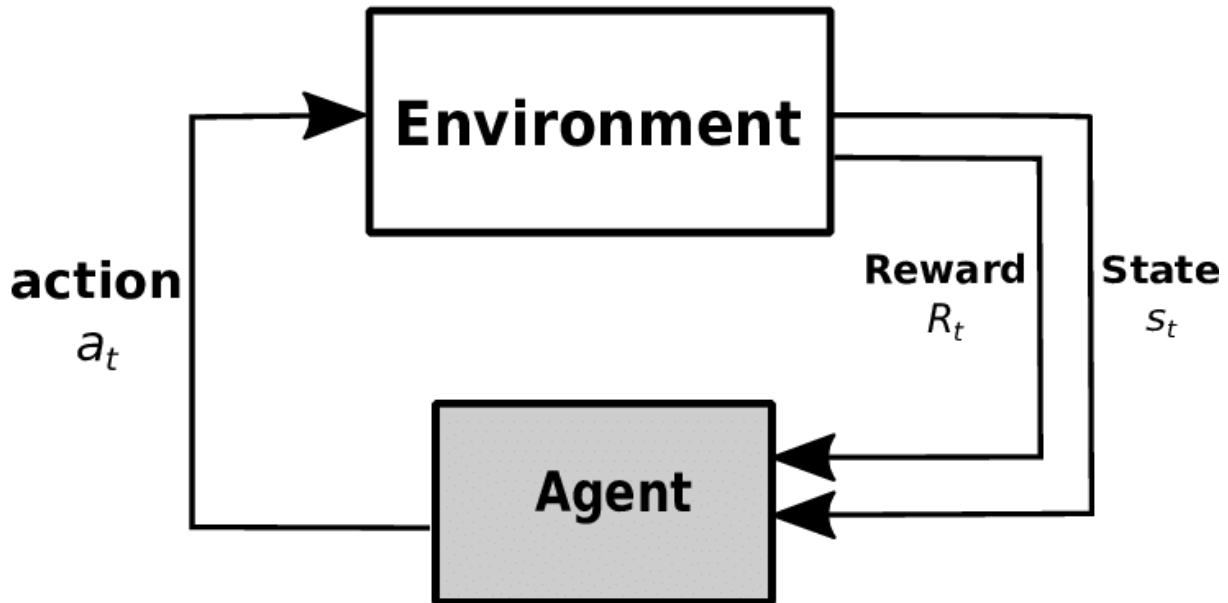
⁷ Regression

⁸ Classification

⁹ Clustering

یادگیری تقویت شده

چنانچه گفته شد، در این پارادایم، معمولا داده ای برای آموزش نداریم (چه به صورت برچسب دار یا بدون برچسب). در عوض دو عنصر مهم با نام های عامل^{۱۰} و محیط^{۱۱} داریم که عامل در اثر تعامل با محیط اطلاعات لازم برای آموزش را جمع آوری میکند (اصطلاحا تجربه کسب میکند). در این پارادایم، ماشین یادگیر سعی دارد که بهترین روش^{۱۲} تعامل با محیط را بیابد.



با توجه به نمودار بالا، ماشین یادگیر بر اساس سیگنالهای جایزه و وضعیت، عملیات a_t را انجام میدهد (که این عملیات میتواند از بین تعداد محدودی عملیات انتخاب شود) و جایزه^{۱۳} و وضعیت^{۱۴} (یا مشاهده^{۱۵} که میتواند گسسته یا پیوسته باشد) جدید را از طرف محیط دریافت میکند در ابتدای کار، عامل یادگیر، ممکن است محیط را کاملا بشناسد و توصیف یا مدلی از آن را داشته باشد (این توصیف مثلا میتواند به صورت یک مدل مارکوف باشد). در این حالت اصطلاحا محیط را کاملا قابل مشاهده^{۱۶} میگویند. عامل ممکن است تنها بخشی^{۱۷} از محیط را بشناسد و یا اساسا هیچ شناخت و مدلی از محیط نداشته باشد^{۱۸}.

چنانچه مدل محیط در دسترس باشد، برای یافتن سیاست بهینه^{۱۹} میتوان از روشهای برنامه نویسی پویا، مونته کارلو، TD و.. استفاده کرده و در صورتی که مدل محیط در دسترس نباشد از الگوریتم هایی نظیر SARSA و Q-Learning استفاده میشود

در این پروژه ما با یک مسئله بدون مدل^{۲۰} مواجهیم و میخواهیم با استفاده از یک الگوریتم بهبود یافته برپلیه Q-Learning، سیاست بهینه را به عامل یادگیر آموزش دهیم لذا در حین تشریح پروژه این الگوریتم و ورژن های بهبود یافته آن را بررسی میکنیم.

¹⁰ agent

¹¹ environment

¹² Optimal policy

¹³ reward

¹⁴ state

¹⁵ observation

¹⁶ Fully observable

¹⁷ Partially observable

¹⁸ Model free

¹⁹ Optimal policy

²⁰ Model free

Q-function:

چنانچه گفته شد هدف ما در این مسئله اینست که عامل یادگیر باتوجه به وضعیتی که دارد (یعنی موقعیت و سرعتش) عملیاتی را انجام دهد (یکی از سه عملیات تعریف شده) که عامل را به هدفش (که خارج شدن از گوال است) نزدیکتر کند. بنابراین به معیاری برای توصیف کارایی و بازده عملیات انتخاب شده نیازمیدیم. این معیار را با تابع $Q(s, a)$ نشان میدهم این تابع وضعیت عامل یادگیر (s) و نیز عملیاتی که در آن وضعیت انتخاب کرده است (a) را دریافت کرده و عددی را به نشانه میزان موثر بودن این عملیات در جهت رسیدن به هدف بازمیگرداند. حال با توجه به این معیار میتوان عملیاتها را مقایسه کرد. اما این معیار چگونه باید تعریف شود؟

قطعا عملیاتی که عامل در هر وضعیت انجام میدهد را باید بر اساس اینکه در نهایت، عامل به هدف میرسد یا نه، مورد سنجش و ارزیابی قرار گیرد. به عبارتی، در ارزیابی عملیات انجام شده در هر وضعیت باید به آینده نیز توجه داشت بر این اساس تابع Q به صورت زیر تعریف میشود:

$$Q(s, a) = E_{\pi}\{R_{t+1} + R_{t+2} + R_{t+3} + \dots | S_t = s, A_t = a\}$$

که یعنی اگر در وضعیت s عملیات a انجام شود و پس از آن نیز عملیات ها بر اساس سیاست π انتخاب شود در نهایت (یا پس از طی تعداد مشخصی $state$) مقدار مورد انتظار مجموع جایزه های دریافتی چقدر است

منظور از **Q-Learning** پیدا کردن تابع Q است!

Q-Learning:

گفتیم که تابع Q را نداریم و لذا باید آنرا تخمین بزنیم. همچنین برای ذخیره کردن این تابع نیز باید راه حلی داشته باشیم. مثلا اگر مسئله ساده باشد میتوان مقادیر Q را در جدول ذخیره کرد یا اگر برای آن شکل پارامتری در نظر بگیریم و پارامترهای آنرا تخمین بزنیم، ذخیره این پارامترها کافیست

Q-Learning الگوریتمیست که به منظور تخمین و یافتن تابع Q استفاده میشود. در این الگوریتم آپدیت کردن مقادیر Q به صورت

زیر انجام میپذیرد:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R + \gamma \max_a Q(s', a) - Q(s, a)]$$

در رابطه بالا α : نرخ بیدادگیری γ : ضریب *discount*

R : جایزه ایست که عامل یادگیر با انجام عملیات a وقتی در وضعیت s قرار دارد دریافت میکند

$\gamma * \max_a Q(s', a)$: بیانگر حداکثر مجموع جایزه قابل انتظار است وقتی عامل از وضعیت s' تا وضعیت نهایی پیش رود. بنابراین جمع این دو ترم به نوعی بیانگر حد اکثر جایزه مورد انتظار است که عامل یادگیر میتواند با عبور از وضعیت جاری (s) به دست آورد. از طرفی تابع Q ، جایزه مورد انتظار را برای عملیات a در وضعیت s مشخص میکند. بنابراین اگر مقادیر Q به عبارت $R + \gamma * \max_a Q(s', a)$ میل کند (از اینرو به ترم مذکور هدف یا **TD target** میگویند). Q حاصله را میتوان بهترین تابع Q ممکن نامید.

Greedy action

حال با فرض داشتن تابع Q ، عامل یادگیر چگونه عملیات متناسب با وضعیت جاری را انتخاب کند؟ فرض میکنیم که روش انتخاب، یک تابع مثلا $\pi(s)$ باشد یعنی این تابع وضعیت s را دریافت کرده و عملیات مناسب در وضعیت s را برگرداند. یک روش آنست که در هر وضعیت، عملیاتی انتخاب شود که بیشترین **Q-value** را داشته باشد در این صورت داریم:

$$\pi(s) = \operatorname{argmax}_a Q(s, a)$$

این روش انتخاب را، حریصانه مینامند. اما اگر عملیاتی که بیشترین Q-value را دارد لزوماً بهترین عملیات نباشد چه؟ پیروی از سیاست حریصانه راه را برای آزمودن سیاستهای دیگر میبندد (در این حالت اصطلاحاً مدل فقط در حال بهره برداریست) در مقابل، عامل یادگیر میتواند به امید یافتن سیاستی بهتر، در فضای عملیاتها اقدام به کاوش کند هرچند، نتیجه این کاوش لزوماً بهتر از سیاست حریصانه نخواهد بود.

ε-Greedy action

در انتخاب عملیات به صورت حریصانه، عملاً احتمال انتخاب دیگر عملیات ها صفر است و لذا کاوشی نداریم. برای اینکه عامل یادگیر امکان کاوش در محیط را داشته باشد باید بخشی از احتمال انتخاب عملیات حریصانه را بین سایر عملیات ها تقسیم کنیم... ! ε-Greedy action به صورت زیر، احتمال را بین عملیاتها تقسیم میکند:

$$\pi(s)_{\epsilon\text{-greedy}} = \begin{cases} \frac{\epsilon}{m} + 1 - \epsilon & \text{if } a = \operatorname{argmax}_a Q(s, a) \\ \frac{\epsilon}{m} & \text{others} \end{cases}$$

یعنی با احتمال $(\frac{\epsilon}{m} + 1 - \epsilon)$ عملیات حریصانه انتخاب میشود و با احتمال $\frac{\epsilon}{m}$ هم بقیه عملیات ها

اما پیاده سازی روش فوق چگونه است؟

پیاده سازی انتخاب با احتمال دلخواه

ایده کلی بدین صورت است که، اگر متغیر تصادفی X دارای توزیع یکنواخت در بازه (0,1) باشد آنگاه داریم

$$P(X > \epsilon) = \int_{\epsilon}^1 1 dx = 1 - \epsilon$$

به عبارت دیگر یک نمونه از این توزیع یکنواخت، با احتمال $1 - \epsilon$ از ϵ بزرگتر است لذا برای پیاده سازی روش ε-Greedy action

به صورت زیر عمل میکنیم:

1. ابتدا یک نمونه (عدد) از توزیع یکنواخت در بازه (0,1) تولید میکنیم
2. اگر نمونه تولید شده از ϵ بیشتر باشد آنگاه عملیات حریصانه را انتخاب کن
3. در غیر این صورت به صورت رندوم یک عملیات را انتخاب کن

Deep Q-Learning (DQN)

در الگوریتم Q-Learning تابع Q را میتوان با یک شبکه عصبی تقریب زد. این حالت را Deep Q-Learning می گویند

overestimation

عامل یادگیر در ابتدای کار هیچ شناختی از محیط ندارد لذا باید با استفاده از اطلاعاتی که در اثر تعامل با محیط جمع آوری میکند فرایند آموزش را پیش ببرد، فقدان شناخت اولیه، منجر به جمع آوری داده هایی میشود که عموماً زائد و حاوی درصد زیادی نویز است. در الگوریتم Q-Learning معمولی، عامل یادگیر به کمک تابع Q و بر اساس روش حریصانه، عملیات بهینه را انتخاب میکند و سپس ارزشیابی و آپدیت کردن تابع Q نیز بر اساس TD target/تفاه می افند ولی TD target خود به Q وابسته است و با تغییر تابع Q تغییر میکند لذا همگرایی به حالت بهینه بسیار دشوار خواهد بود. از طرفی با توجه به نویز بالای داده های جمع آوری شده، عملیاتی که بیشترین Q-value را دارد لزوماً

عملیات بهینه نخواهد بود که مورد آخر در کنار سایر موارد میتواند منجر به پدیده ای موسوم به overestimation شود. به طول خلاصه اگر مدل دچار overestimation شود، در این صورت در هر وضعیت عملیاتی انتخاب میشود که دارای بیشترین Q-value است ولی بهینه نیست! به منظور حل این مشکل راهکارهای مختلفی ارائه شده است. یکی از این راهکارها استفاده از دو تابع Q و Q' یکی برای انتخاب عملیاتها و دیگری برای ارزشیابی عملیات این ایده راه Double Q-Learning گویند...

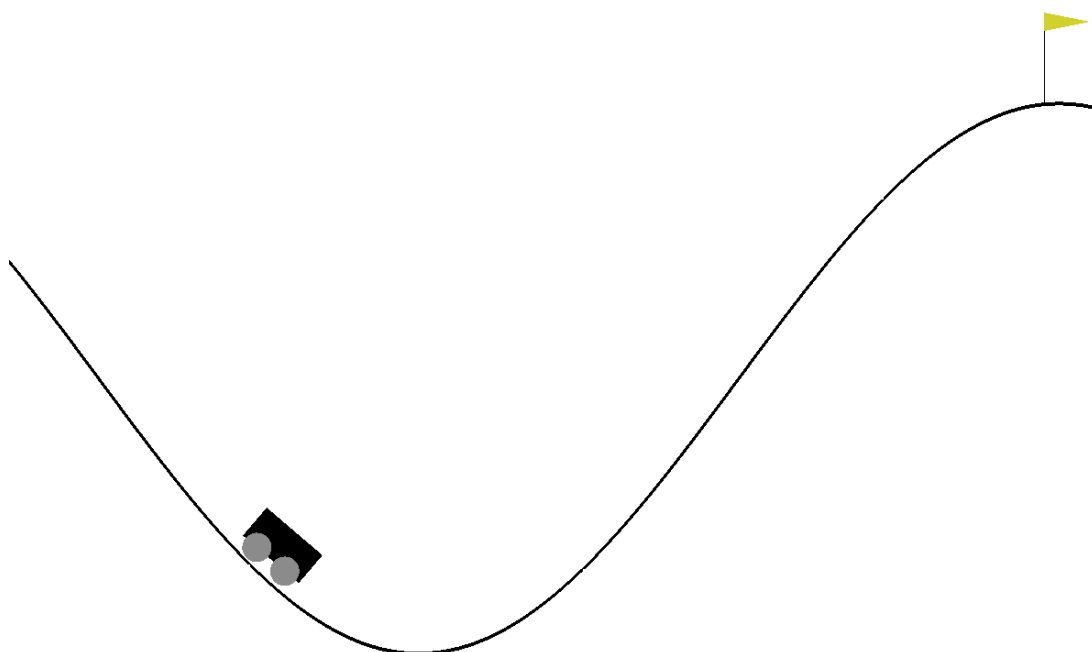
Double DQN(DDQN)

همان Double Q-Learning است که توابع Q و Q' را با شبکه های عصبی تقریب زده میشود. بنابراین در این الگوریتم باید دو تابع Q و Q' را پیدا کنیم که این امر پیچیدگی محاسباتی مسئله را بیشتر میکند لذا به منظور کاهش بار محاسباتی از روشی که در (مقاله پیوست شده پروژه) تشریح شده است استفاده میکنیم در این روش، اگرچه دو تابع داریم (دو شبکه عصبی) ولی فقط یکی از آنها آموزش داده میشود (train_newtwork) و پس از تعداد مشخصی گام آموزشی، وزنهاى شبکه آموزش دیده به شبکه دوم (target_network) منتقل میشود.

بخش دوم

تشریح مسئله

مسئله Mountain Car یکی از اولین و مشهورترین مسائل در حوزه یادگیری تقویتی است. در این مسئله عامل، ماشینی است که در یک گودال (محیط مسئله) گیر افتاده و باید در طی تعامل با محیط خود را به نقطه هدف (نشانه داده شده با پرچم) برساند اما چون نیروی موتور این ماشین برای غلبه بر گرانش کافی نیست نمیتواند با حرکت مستقیم در جهت هدف، خود را به هدف برساند. بنابراین ماشین باید ابتدا در جهت عکس هدف از شیب مقابل بالا رفته و پس از کسب انرژی پتانسیل لازم خود را رها کند و در نهایت در نقطه ای مشخص، با روشن کردن موتور خود را به هدف برساند!... هر اپیزود از این مسئله با انجام 200 عملیات و یا رسیدن به نقطه هدف، پایان میابد



بر اساس کدهای مربوطه به این مسئله در gym، گودال درواقع یک منحنی سینوسی به معادله زیر است:

$$y = 0.45 * \sin(3 * x) + 0.55 \quad -1.2 < x < 0.6$$

و لذا با توجه به رابطه بالا که منحنی مسیر هم نامیده میشود رابطه سرعت باید کسینوسی باشد. در این مسئله، عامل یادگیر میتواند یکی از سه عملیات زیر را انجام دهد:

0. حرکت به چپ، 1. حرکت نکند، 2. حرکت به راست

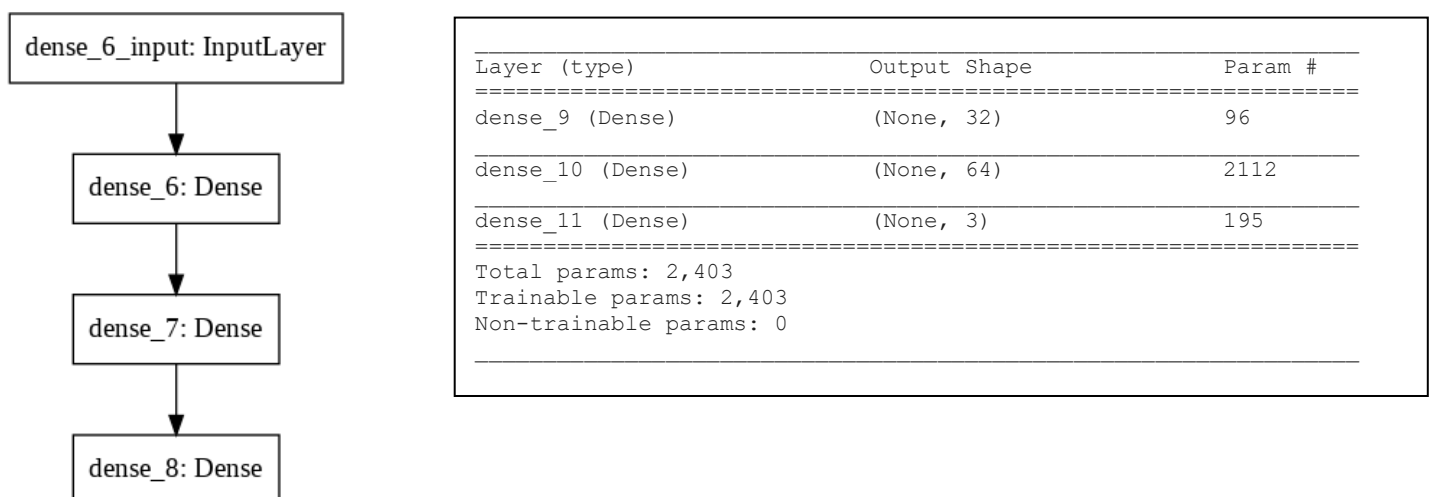
عامل یادگیر پس از انجام هر عملیاتی جایزه 1- دریافت میکند مگر زمانی که عملیات انجام شده منجر به رسیدن به هدف شود که در این حالت جایزه 0 را دریافت میکند همچنین وضعیت (State) نیز به ترتیب شامل دو مولفه مکان عامل و سرعت آن میباشد

$$state = (-1.2 < position < 0.6, 0.07 < velocity < 0.07)$$

هدف از حل مسئله، یافتن دنباله ای از عملیاتی مجاز (عملیات های 0 و 1 و 2) است که دنبال کردن آن منجر به خروج ماشین از گودال شود (بهترین راه حل هم دنباله ایست که منجر به بیشترین جایزه شود). در ادامه با طراحی آموزش یک شبکه عصبی سعی میکنیم Q را برای این مسئله را بیابیم

شبکه های عصبی به عنوان تقریبی از تابع سیاستگذاری [Policy function]

در این بخش دو شبکه عصبی برای تقریب زدن توابع Q و Q' میسازیم. شبکه model_for_train معادل Q و شبکه model_for_predict معادل Q' است. ورودی این شبکه ها زوج مرتب (مکان، سرعت) است که هردو کمیتی پیوسته اند همچنین خروجی شبکه ها نیز یک بردار سه مولفه ای است که مولفه های آن به ترتیب ارزش عملیاتی 0 1 2 را مشخص میکنند. چون ارزش هر عملیات، کمیتی پیوسته است لذا با یک مسئله رگرسیون مواجههیم. ساختار کلی هر دو شبکه یکسان بوده و به صورت زیر است:



چنانچه گفته شده تنها وزنه های شبکه model_for_train آپدیت شده و پس چند مرتبه آپدیت این وزنها به مدل دیگر انتقال میابد برای آپدیت کردن هم از رابطه زیر به عنوان target شبکه استفاده میشود:

$$y^{DDQN} = R_{t+1} + \gamma * model_for_predict(S_{t+1}, \operatorname{argmax}_a model_for_train(S_{t+1}, a))$$

در کادر زیر تعدادی از هایپر پارمترهای مسئله آورده شده است

چون عامل یادگیر در ابتدای کار هیچ شناختی از محیط ندارد لذا بهتر است در اپیزودهای اولیه حداکثر کاوش را داشته باشیم به همین دلیل مقدار اولیه ایسیلون 0.99 انتخاب شده است و در هر اپیزود 0.005 از آن کم میشود تا نهایتاً با 0.01 برسد.

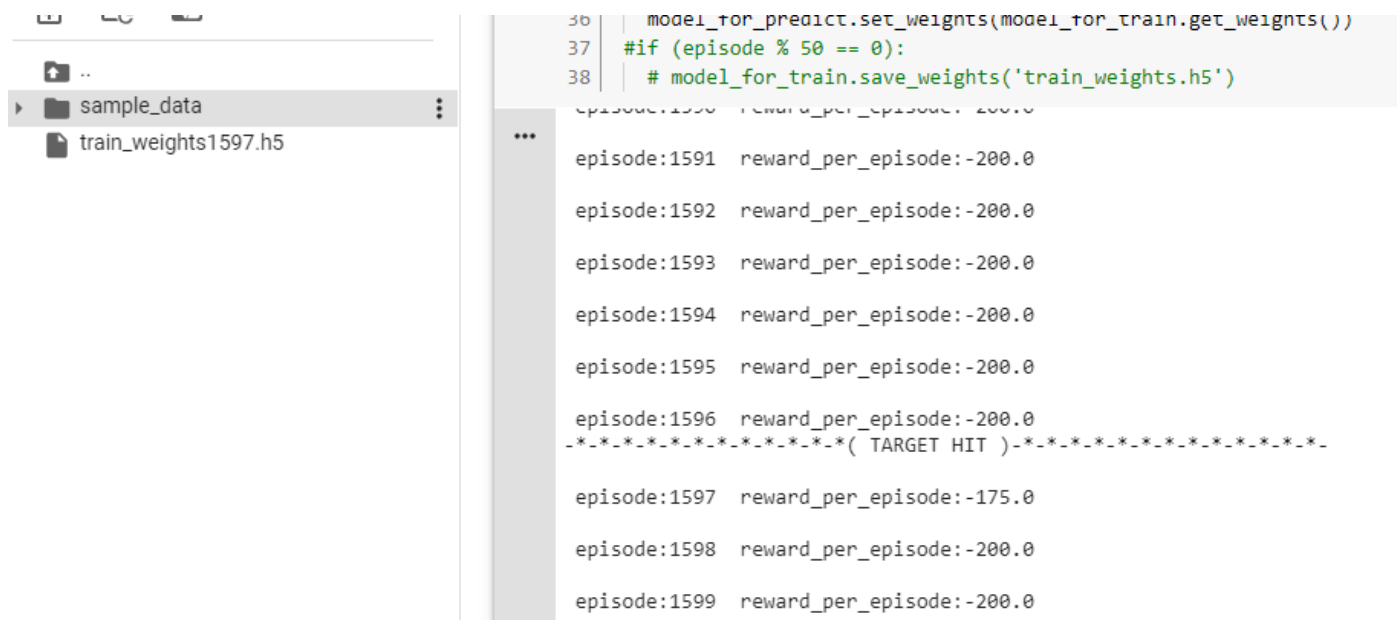
(*) چون امکان استفاده از `env.render()` در گوگل کولب وجود نداشت!، عامل یادگیر در گوگل کولب آموزش داده شده و سپس وزنهای حاصله در سیستم شخصی، آزمایش شده است (فایل `test.py`)

(*) فرایند یادگیری چندین بار و به طور کامل از اول اجرا شده است (یعنی تمام متغیرها در مقادیر پیشفرضشان بوده اند و memory هم خالی بوده است) و هر بار پس از تعداد نامشخصی اپیزود، جواب مسئله به دست آمد. بنابراین میتوان نتیجه گرفت که الگوریتم اسفاده شده نسبت به شرایط اولیه بسیار حساس است.

(*) در اولین اجرا، در اپیزود 1301 به جواب رسیدیم (عبارت TARGET HIT قبل از اعلام شماره اپیزوده پرینت میشود) و پس از آن در چندین اپیزود بعدش هم، عامل یادگیر به هدف رسیده است. اما با ادامه فرایند یادگیری، وضعیت بهتر نمیشود.

```
episode:1297 reward_per_episode:-200.0  
...  
episode:1298 reward_per_episode:-200.0  
  
episode:1299 reward_per_episode:-200.0  
  
episode:1300 reward_per_episode:-200.0  
*_**_*_*_*_*_*_*_*_*( TARGET HIT )-*_*_*_*_*_*_*_*_*_*_*_  
  
episode:1301 reward_per_episode:-124.0  
*_**_*_*_*_*_*_*_*_*( TARGET HIT )-*_*_*_*_*_*_*_*_*_*_*_  
  
episode:1302 reward_per_episode:-175.0  
*_**_*_*_*_*_*_*_*_*( TARGET HIT )-*_*_*_*_*_*_*_*_*_*_*_  
  
episode:1303 reward_per_episode:-176.0  
  
episode:1304 reward per episode:-200.0
```

در اجرای دوم، پس از 1597 ایزود هدف مشاهده شد. (و وزنهای آن ذخیره شد)



```
episode:309 reward_per_episode:-190.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:310 reward_per_episode:-150.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:311 reward_per_episode:-121.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:312 reward_per_episode:-179.0  
  
episode:313 reward_per_episode:-200.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:314 reward_per_episode:-154.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:315 reward_per_episode:-157.0  
-_*_-*_-*_-_*_-_*_-*( TARGET HIT )-_*_-*_-*_-_*_-_*_-_*_-_*_-_*_-_*_-  
  
episode:316 reward per episode:-162.0
```

```
..
train_weights21.h5
train_weights22.h5
train_weights298.h5
train_weights300.h5
train_weights301.h5
train_weights302.h5
train_weights303.h5
train_weights304.h5
train_weights305.h5
train_weights306.h5
train_weights307.h5
train_weights309.h5
train_weights310.h5
train_weights311.h5
train_weights312.h5
train_weights314.h5

30 | break
31 |
32 | current_state = new_state.reshape(1,2)
33 |
34 | print("\n episode:{} ".format(episode)+" reward_per_episode:{}".format(reward_per_episode))
35 | train_on_batch(32)
36 | if(episode %2 ==0):
37 |     model_for_predict.set_weights(model_for_train.get_weights())
38 | #if (episode % 50 == 0):
39 |     # model_for_train.save_weights('train_weights_{}.h5'.format(episode))

...
episode:357 reward_per_episode:-200.0
episode:358 reward_per_episode:-200.0

[ ] 1 model_for_train.save_weights("train_weights_{}.h5".format(episode))
```

یعنی در همان اپیزودهای ابتدایی (اپیزود چهار تا بیست و بست و دو) و سپس در اپیزودهای بالاتر از 300 به طور مداوم عامل به هدف میرسید و وزنهای مربوطه هم ثابت شده است (هرچند باز هم با ادامه فرایند، عامل از هدف دور شد). دلیل این امر احتمالاً خالی شدن memory باشد و به تعبیری فراموشی مدل است!...در ادامه به بررسی این اثر میپردازیم

نقش مکانیزم فراموشی :

عامل یادگیر نتیجه تعاملات خود با محیط را در حافظه خود (متغیر memory) ذخیره میکند اما چنانچه گفته شد بخش قابل توجهی از این اطلاعات ذخیره شده زائد و همراه با نویز زیاد است لذا احتمالاً فراموش کردن برخی از این اطلاعات منجر به بهبود وضعیت یادگیری شود. به منظور ایجاد مکانیزمی برای فراموشی از کد زیر در فرایند یادگیری استفاده شد. این کد، حافظه را پس از هر 512 داده ثبت شده پاک میکند

```
if(len(memory) == 512):
    memory.clear()
```

در اولین اجرا نتیجه زیر حاصل شد:

در این پروژه با برخی از الگوریتم‌های مطرح در یادگیری تقویتی آشنا شدیم و عوامل تاثیر گذار بر این الگوریتم‌ها را شناختیم و همچنین با کتابخانه *gym* که ابزاری آماده برای شبیه سازی الگوریتم‌های حوزه یادگیری تقویتی است نیز آشنا شدیم. سپس یکی از این الگوریتم‌ها را پیاده سازی کرده و مشکلات پیاده سازی عملی آنها را بررسی نمودیم