

دانشگاه صنعتی خواجه نصیرالدین طوسی  
دانشکده مهندسی برق - گروه مهندسی کنترل

## درس یادگیری ماشین مینی پروژه سوم

نام و نام خانوادگی	علی مهربابی
شماره دانشجویی	۴۰۲۲۳۸۵۴
تاریخ	اردیبهشت ۱۴۰۳
استاد درس	دکتر علیاری



## فهرست مطالب

۴	۱	لینک های مربوطه
۴	۱.۱	لینک مربوط به گیت هاب
۴	۲.۱	لینک مربوط به گوگل کلب
۵	۲	سوال دوم
۵	۱.۲	قسمت اول
۶	۲.۲	قسمت دوم
۱۶	۳.۲	قسمت سوم



## فهرست تصاویر

۵	.....	محیط Lunar Lander	۱
۷	.....	import Lunar Lander from Gym	۲
۱۳	.....	batch size = 32 نمودار پاداش تجمعی برای	۳
۱۳	.....	batch size = 64 نمودار پاداش تجمعی برای	۴
۱۴	.....	batch size = 128 نمودار پاداش تجمعی برای	۵
۱۴	.....	batch size = 32 پاداش تجمعی برای	۶
۱۵	.....	batch size = 64 پاداش تجمعی برای	۷
۱۵	.....	batch size = 128 پاداش تجمعی برای	۸
۱۶	.....	batch size = 128 ویدیویی از ایجنت در محیط به ازای در اپیزود ۲۵۰	۹
۱۶	.....	batch size = 128 ساختن ویدیو از ایجنت در محیط به ازای در اپیزود ۲۵۰	۱۰
۱۷	.....	batch size = 128 نمودار پاداش تجمعی به کمک DDQN در حالت	۱۱
۱۸	.....	batch size = 128 ویدیویی از ایجنت در محیط به ازای در اپیزود ۲۵۰	۱۲



## فهرست برنامه‌ها

۶	.....	swig installing	۱
۷	.....	Xvbf installing	۲
۷	.....	GPU using	۳
۸	.....	display virtual a initializing	۴
۸	.....	replay Experience	۵
۹	.....	class DeepQNetwork	۶
۱۰	.....	action take	۷
۱۱	.....	params update	۸
۱۱	.....	weights the load and save	۹
۱۱	.....	۳۲ size batch	۱۰
۱۲	.....	agent initialize	۱۱
۱۲	.....	loop train	۱۲
۱۲	.....	loop episode	۱۳
۱۲	.....	Decay Epsilon	۱۴
۱۷	.....	Agent DQN	۱۵
۱۷	.....	Agent DDQN	۱۶



## ۱ لینک های مربوطه

### ۱.۱ لینک مربوط به گیت هاب

از این لینک **گیت هاب (Github)** می توانید برای دسترسی به صفحه Github مربوط به این پروژه استفاده کنید.

### ۲.۱ لینک مربوط به گوگل کلب

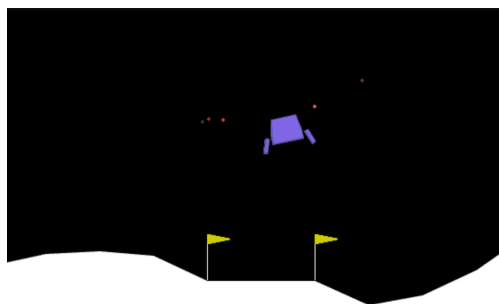
از این لینک **گوگل کلب (Google Colab)** می توانید برای دسترسی به notebook نوشته شده دسترسی پیدا کنید.



## ۲ سوال دوم

## ۱.۲ قسمت اول

محیط Lunar Lander یک محیط توسعه یافته توسط Open AI در قالب یک کتابخانه متن باز به نام Gym است. از این قالب برای شبیه سازی ها و اهداف آموزشی برای آموزش Reinforcement Learning استفاده می شود. این کتابخانه شامل محیط های زیادی است که هریک ما را با جنبه ای از RL آشنا می کند. مسائل برگرفته از کنترل کلاسیک مانند Cart Pole تا بازی آتاری و حتی شبیه سازی های پیچیده تر مانند MuJoCo. از خوبی های این منبع متن باز؛ یکپارچه بودن آن در تمامی فضاهاست به طوری که تمامی فضاها و نتایج مربوط به آن با هم قابل مقایسه هستند. متد هایی مانند `close()`، `reset()` یا `step(action)` از جمله این متد ها هستند. حالا که با فضای کلی، توسعه دهنده و اهداف این کتابخانه متن باز آشنا شدیم به سراغ یکی از محیط های آن یعنی Lunar Lander می رویم. شکل ۱ تصویری از این محیط نشان می دهد که در آن یک فضاپیما قصد دارد روی یک مساحت مشخص فرود بیاید. در واقع این یک مسئله بهینه سازی مسیر این راکت یا فضاپیما است که یک مساله کلاسیک است. با توجه به اصل Pontryagin's maximum حالتی بهینه است که موتور در حالت حداکثری خود روشن باشد یا خاموش شود. به همین علت هر موتور دو حالت دارد؛ روشن باشد یا خاموش باشد.



شکل ۱: محیط Lunar Lander

Action Space در این بازی شامل چهار اکشن گسسته است:

- کاری انجام نشود
- موتور جهت چپ روشن شود
- موتور اصلی روشن شود
- موتور جهت راست روشن شود

State Space در این بازی شامل یک بردار ۸ بعدی است.

- پوزیشن  $X$  که می تواند بین  $-1.5$  تا  $1.5$  باشد.
- پوزیشن  $Y$  که می تواند بین  $-1.5$  تا  $1.5$  باشد.



- سرعت خطی در جهت  $X$  که می تواند بین  $-5$  تا  $5$  باشد.
- سرعت خطی در جهت  $Y$  که می تواند بین  $-5$  تا  $5$  باشد.
- در جه که می تواند بین  $-3.14$  تا  $3.14$  (به رادیان) باشد.
- سرعت زاویه ای که می تواند بین  $-5$  تا  $5$  باشد.
- یک متغیر بولین که نشان می دهد پای ۱ با زمین ارتباط دارد یا نه
- یک متغیر بولین که نشان می دهد پای ۲ با زمین ارتباط دارد یا نه

Reward System به این صورت است که بازیگر یا فضاییما را ترغیب می کند تا به صورت دقیق و نرم در محل موردنظر فرود بیاید. پاداش برای حرکت از بالا و فرود آمدن روی سکوی مورد نظر ۱۰۰ تا ۱۴۰ امتیاز است. اگر فضاییما از سکویی که باید روی آن فرود بیاید فاصله بگیرد پاداش از دست می دهد. اگر تصادف بکند -۱۰۰ امتیاز از دست می دهد. اگر روی سکو بنشیند و بی حرکت بماند یا rest انجام دهد ۱۰۰ امتیاز دریافت می کند. اگر هر پا با زمین برخورد کند ۱۰ امتیاز می گیرد. راه افتادن موتور اصلی  $-0.03$  امتیاز در هر فریم را در بر دارد و برای موتورهای کناری نیز همین مقدار است. اینکار برای تشویق به استفاده کمتر از سوخت انجام می گیرد. رسیدن به حالت solved ۲۰۰ امتیاز دارد. پوزیشن ابتدایی فضاییما در مرکز و بالای تصویر است و در ابتدا یک نیروی رندوم به مرکز جرم آن اعمال می شود. هر اپیزود بنابر یکی از دلایل زیر تمام می شود:

- فضاییما تصادف کند. یعنی بدن آن با ماه برخورد کند.
- فضاییما خارج از فضای دید بشود برای مثال حرکت افقی آن بیشتر از ۱ شود.
- اگر فضاییما بیدار نباشد، یعنی حرکتی نداشته باشد و با چیز دیگری برخورد نکند.

## ۲.۲ قسمت دوم

برای حل این قسمت از کد قرار گرفته به عنوان نمونه استفاده شده و تغییراتی در آن اعمال شده تا ارورهای آن برطرف شود. در مرحله اول ابتدا باید موارد و کتابخانه های لازم را نصب کنیم. دو نکته وجود دارد که قابل اشاره هستند، اول اینکه باید دو دستور برای نگرفتن ارور به کد اضافه کنیم که در برنامه ۱ آمده اند.

```
1 !apt-get update
2 !apt-get install -y swig
```

Code 1: installing swig

دوم برای نشان دادن ویدیوها به صورت مجازی باید از Xvfb استفاده کنیم که نیاز است تا آن را در محیط google collab نصب کنیم که به وسیله برنامه ۲ اینکار انجام شده است. سایر دستورات در بخش نصب ثابت باقی مانده است و به سراغ مرحله بعد می رویم. در مرحله بعد دستور مورد نیاز برای استفاده از GPU در محیط google collab را به صورت برنامه ۳ وارد می کنیم که در نتیجه به ما می گوید که ما در حال استفاده از Cuda و در حال ران کردن GPU هستیم.



```
1 !sudo apt-get update
2 !sudo apt-get install xvfb
```

Code 2: installing Xvfb

```
1 import torch
2 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
3 device
```

Code 3: using GPU

در ادامه به صورت **شکل ۲** محیط مورد نظر که همان Lunar Lander است را از Gym وارد می کنیم. همانطور که مشاهده می شود ما ۸ observation یا state داریم و می توانیم ۴ action انجام دهیم که همگی در **بخش ۱.۲** توضیح داده شد. در مرحله بعد ابتدا یک صفحه نمایش مجازی تعریف می کنیم و سائز آن را مشخص می کنیم و سپس یک

```
[5] # enviroment
import gym
env = gym.make('LunarLander-v2', render_mode="rgb_array")
#TODO: find observation size: 8
state_size = env.observation_space.shape[0]
#TODO: find action size: 4: 0- Do nothing 1- Fire left engine 2- Fire down engine 3- Fire right engine
action_size = env.action_space.n
state_size, action_size
```

(8, 4)

شکل ۲: محیط import Lunar Lander from Gym

تابع تعریف می کنیم تا بتوانیم ویدیو هارا که در آینده قصد ضبط آن هارا داریم در محیط Jupiter notebook ببینیم. انجام اینکار در **برنامه ۴** آمده است.

در مرحله بعد **برنامه ۵** را داریم. در ابتدای آن یک تاپل می سازیم که حالت یک experience را در خود ذخیره کند. این موارد شامل حالتی که فضاییما در آن وجود دارد (قبل از انجام اکشن)، اکشنی که انجام می دهد، حالت بعدی که با انجام اکشن به آن می رود، پاداشی که دریافت می کند و یک متغیر بولین که نشان می دهد آیا عمل تمام شده یا نه است. سپس کلاس ExperienceReplay یک reply buffer را پیاده می کند. در واقع تمامی experience های گذشته یا همان transitions را در خود ذخیره می کند تا در مرحله آموزش از آنها استفاده شود. برای اینکار از یک deque یا double-ended queue استفاده می شود تا memory ساخته شود. reply buffer به agent اجازه می دهد تا از تجربیات گذشته خود با سمپل کردن به صورت رندوم از batch های transition یاد بگیرد و این رندوم سازی باید پایداری یادگیری می شود و موضوع correlation بین experience های متوالی را برای ما بهبود می دهد. store trans ترنژیشن های جدید را به reply buffer اضافه می کند.

در مرحله بعد شبکه عصبی عمیق مربوط به الگوریتم DQN را تشکیل می دهیم. **برنامه ۶** این شبکه را نشان می دهد که در کلاس DeepQNetwork تعریف شده است. همانطور که مشخص است این یک شبکه sequential است که نسبت به کد اصلی به صورت مستقیم پیاده سازی شده و برگردانده می شود. ورودی این شبکه به اندازه state ها نرون دارد و در نهایت خروجی این شبکه همان Q value ها هستند که مقادیر انتظاری از پاداش آینده را برای هر اکشن به ما می دهد. پس خروجی این شبکه به اندازه اکشن ها است و مقادیر Q برای آن ها تخمین زده می شود. در لایه اول یک





```

1 display = Display(visible=0, size=(1400, 900))
2 display.start()
3
4 def show_video():
5     mp4list = glob.glob('video/*.mp4')
6     if len(mp4list) > 0:
7         mp4 = mp4list[0]
8         video = io.open(mp4, 'r+b').read()
9         encoded = base64.b64encode(video)
10        ipythondisplay.display(HTML(data='''<video alt="test" autoplay loop controls style="height:
    ↪ 400px;">
11                <source src="data:video/mp4;base64,{0}" type="video/mp4" />
12                </video>'''.format(encoded.decode('ascii'))))
13    else:
14        print("Could not find video")

```

Code 4: initializing a virtual display

```

1 Transition = namedtuple('Transition', ('state', 'action', 'next_state', 'reward', 'done'))
2
3 class ExperienceReplay():
4     def __init__(self, capacity):
5         self.memory = deque([], maxlen=capacity)
6
7     def store_trans(self, s, a, sp, r, done):
8         transition = Transition(s, a, sp, r, done)
9         self.memory.append(transition)
10
11    def sample(self, batch_size):
12        return random.sample(self.memory, batch_size)
13
14    def __len__(self):
15        return len(self.memory)

```

Code 5: Experience replay



لایه تماماً متصل با تابع فعال ساز ReLU استفاده شده است. سپس از Layer normalization استفاده شده تا روند یادگیری را پایدار سازد و سرعت بیخشد، همچنین برای جلوگیری از overfitting از روش Dropout استفاده شده است و نرخ آن ۱۰ درصد است. لایه دوم نیز یک لایه خطی به صورت تماماً متصل است و دوباره از نرمالیزه کردن و Dropout استفاده شده و تابع فعال ساز ReLU است. لایه سوم نیز به همین صورت. متد forward ورودی تانسور  $x$  را به شبکه می دهد و در خروجی Q value های مربوط به هر اکشن بازگردانده می شود.

```

1 class DeepQNetwork(nn.Module):
2     def __init__(self, state_size, action_size):
3         super(DeepQNetwork, self).__init__()
4         self.net = nn.Sequential(
5             nn.Linear(state_size, 512),
6             nn.ReLU(),
7             nn.LayerNorm(512),
8             nn.Dropout(0.1),
9             nn.Linear(512, 512),
10            nn.ReLU(),
11            nn.LayerNorm(512),
12            nn.Dropout(0.1),
13            nn.Linear(512, 512),
14            nn.ReLU(),
15            nn.Linear(512, action_size)
16        )
17
18    def forward(self, x):
19        return self.net(x)

```

Code 6: DeepQNetwork class

کلاس بعدی DQNAgent است که Agent ما را پیاده سازی می کند. ایجنتی که با محیط interact می کند، experience ها را ذخیره می کند و Q value ها را بروز رسانی می کند. ابتدا نیاز است تا متغیر ها را در کلاس initialize کنیم. این موارد شامل ابعاد state و action و یا batch size است. همچنین نیاز به تعریف gamma داریم که همان discount factor است. tau پارامتر soft update برای شبکه هدف است. همچنین نیاز به experience replay و یک بهینه ساز و چند مورد دیگر نیز داریم.

متد بعدی در این کلاس take action است که در برنامه ۷ آمده است. که بر اساس epsilon greedy policy از حالت فعلی یک اکشن انتخاب می کند. اگر یک عدد رندوم از eps بزرگتر باشد، ایجنت اکشنی را انتخاب می کند که بیشترین Q-value که در شبکه اصلی پیش بینی شده دارد. در غیر این صورت ایجنت یک اکشن رندوم انتخاب می کند و این تفاوت بین exploitation و exploration است. این روش مطمئن می شود که بالانس بین این دو روش در آموزش رعایت می شود.

متد بعدی update params است که در برنامه ۸ آمده است. این متد پارامتر های شبکه اصلی را به کمک تجربه های گرفته شده از replay buffer بروز رسانی می کند و یک soft update را روی شبکه هدف انجام می دهد. در واقع به کمک sample(batch size) یک بیج از تجربه را از replay buffer بر می دارم و با فرمول بلمن مقدار Q value را حساب



```

1  def take_action(self, state, eps=0.0):
2      self.value_net.eval()
3      if random.random() > eps:
4          with torch.no_grad():
5              return torch.argmax(self.value_net(torch.tensor(state).float().unsqueeze(0).to(
6  ↪ device))).item()
7      else:
8          return np.random.randint(0, self.action_size)

```

Code 7: take action

میکند. سپس اختلاف بین مقدار پیش بینی شده و اصلی را می یابد و یک گرادین کاهشی برای مینیمایز کردن این خطا اجرا می کند.

در نهایت به کمک برنامه ۹ این وزن هارا ذخیره و بارگزاری می کند. حالا باید برای سه batch مختلف، عمل آموزش را انجام دهیم. توجه شود که تمامی مراحل توضیح داده شده و همچنین آموزش برای سه حالت یکسان است و فقط batch size تغییر می کند. برای مثال برای بچ ۳۲، به صورت برنامه ۱۰ خواهد بود.

حال به توضیح فرآیند آموزش می پردازیم. ابتدا به صورت برنامه ۱۱ ایجنت را تعریف می کنیم و همچنین پاداش تجمعی را initialize می کنیم. در این قسمت ما یک ایجنت از DQN Agent با حالت و اکشن خاص درست می کنیم که بچ سایز اش را مشخص کرده ایم. همچنین صف های ۲۵ تایی تشکیل می دهیم که در آینده از آنها میانگین می گیریم.

در برنامه ۱۲ حلقه آموزش را به تعداد مشخصی از اپیزود که در این حالت ۲۵۰ است تشکیل می دهیم. در این کد از هر ۵۰ امین اپیزود فیلم برداری می شود یعنی اپیزود ۵۰ ام ۱۰۰ ام ۱۵۰ ام ۲۰۰ ام و ۲۵۰ ام. سپس محیط را ریست می کند و مقادیر اولیه را میگیرد. سپس done را که قبل تر توضیح داده شده بود به حالت false در می آورد تا اپیزود بعدی استارت شود. سپس مقدار پاداش برای اپیزود فعلی را صفر می گذارد.

بخش برنامه ۱۳ تا زمانی که اپیزود تمام شود با محیط ارتباط می گیرد، تجربه کسب می کند و Q value هارا به روز می کند. در واقع ابتدا یک اکشن انتخاب می شود، حالت بعدی، فلگ و پاداش را میگیرد و تجربه کسب شده را در reply buffer ذخیره می کند. سپس Q value function را به روز رسانی می کند (به کمک سمپل گرفته از reply buffer) سپس حالت را از حالت فعلی به حالت بعدی تغییر میدهد و پاداش را به پاداش های قبلی در آن اپیزود اضافه میکند.

بخش بعدی اپسیلون را کاهش می دهد و مقدار پاداش تجمعی در این اپیزود را ذخیره می کند. اگر اپیزود مضربی از ۵۰ بود آن را ذخیره می کند و اگر مضرب ۲۵ بود، پاداش میانگین آم و اپیزود فعلی و مقدار اپسیلون را پرینت می کند. این بخش در برنامه ۱۴ آمده است.

در نهایت پاداش تجمعی برای batch size = 32 در شکل ۳ آمده است. همین نمودار برای batch size = 64 در شکل ۴ آمده است. که همانگونه که دیده میشود وضعیت پایدار تری را دارد ولی در انتها پاداش batch size = 32 بهتر است و نشان می دهد ایجنت بهتر عمل کرده است.

همین نمودار برای batch size = 128 در شکل ۵ آمده است. که همانگونه که دیده میشود وضعیت ناپایدار تری را دارد. مطابق شکل ۳ تا شکل ۵ در مراحل اولیه مقدار پاداش نوسان زیادی دارد که نشان دهنده فاز exploration ایجنت



```

1 def update_params(self):
2     if len(self.experience_replay) < self.batch_size:
3         return
4     batch = Transition(*zip(*self.experience_replay.sample(self.batch_size)))
5
6     state_batch = torch.tensor(np.array(batch.state), dtype=torch.float32).to(device)
7     action_batch = torch.tensor(batch.action, dtype=torch.int64).unsqueeze(1).to(device)
8     next_state_batch = torch.tensor(np.array(batch.next_state), dtype=torch.float32).to(
9         ↪ device)
10    reward_batch = torch.tensor(batch.reward, dtype=torch.float32).unsqueeze(1).to(device)
11    done_batch = torch.tensor(batch.done, dtype=torch.float32).unsqueeze(1).to(device)
12
13    q_expected = self.value_net(state_batch).gather(1, action_batch)
14    q_targets_next = self.target_net(next_state_batch).detach().max(1)[0].unsqueeze(1)
15    q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch))
16    loss = F.mse_loss(q_expected, q_targets)
17
18    self.optimizer.zero_grad()
19    loss.backward()
20    self.optimizer.step()
21
22    # Soft update target network parameters
23    for target_param, local_param in zip(self.target_net.parameters(), self.value_net.
    ↪ parameters()):
24        target_param.data.copy_(self.tau * local_param.data + (1.0 - self.tau) * target_param
    ↪ .data)

```

Code 8: update params

```

1 def save(self, fname):
2     torch.save(self.value_net.state_dict(), fname)
3
4 def load(self, fname):
5     self.value_net.load_state_dict(torch.load(fname, map_location=device))

```

Code 9: save and load the weights

```

1 # NOTE: DON'T change values
2 n_episodes = 250
3 eps = 1.0
4 eps_decay_rate = 0.97
5 eps_end = 0.01
6 BATCH_SIZE = 32

```

Code 10: batch size 32



```

1 agent = DQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
2 crs = np.zeros(n_episodes)
3 crs_recent = deque(maxlen=25)

```

Code 11: initialize agent

```

1 for i_episode in range(1, n_episodes + 1):
2     env = RecordVideo(gym.make("LunarLander-v2", render_mode="rgb_array"), f"./DQN/batch{
    ↪ BATCH_SIZE}/eps{i_episode}") if i_episode % 50 == 0 else gym.make("LunarLander-v2",
    ↪ render_mode="rgb_array")
3     state, info = env.reset()
4     done = False
5     cr = 0

```

Code 12: train loop

```

1 while not done:
2     action = agent.take_action(state, eps)
3     next_state, reward, done, truncated, info = env.step(action)
4     agent.experience_replay.store_trans(state, action, next_state, reward, done or truncated)
5     agent.update_params()
6     state = next_state
7     cr += reward

```

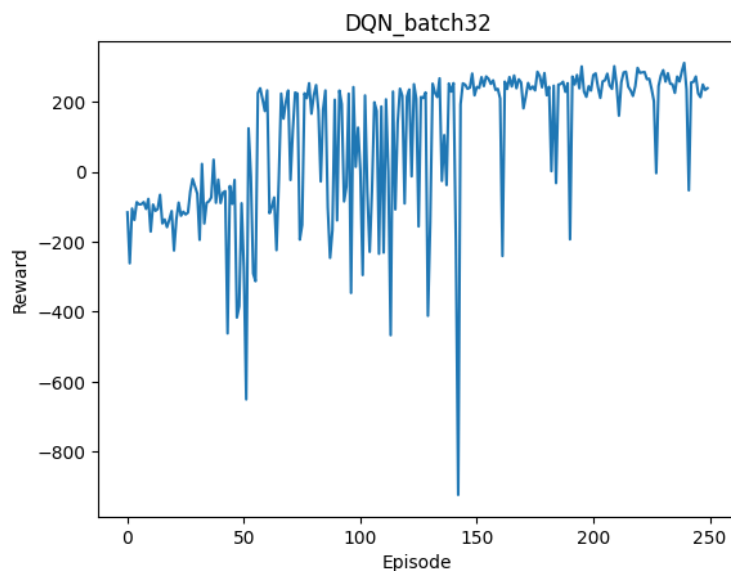
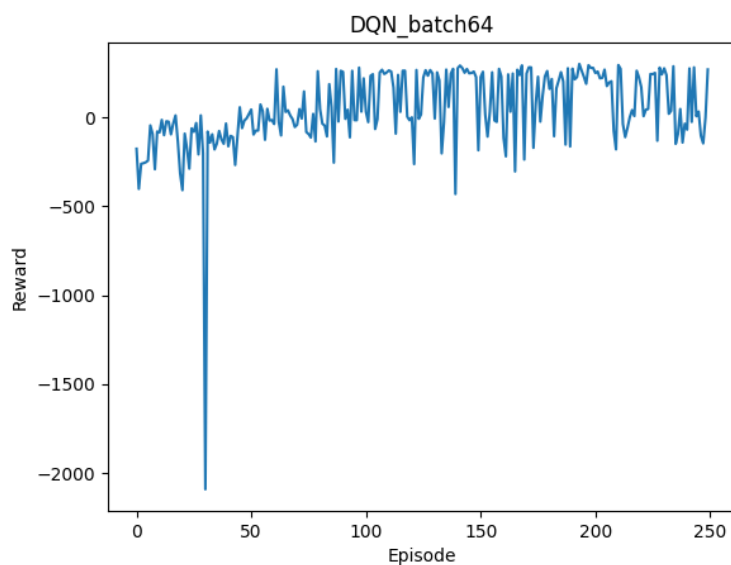
Code 13: episode loop

```

1 eps = max(eps * eps_decay_rate, eps_end)
2 crs[i_episode - 1] = cr
3 crs_recent.append(cr)
4 if i_episode % 50 == 0:
5     agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
6
7 print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps:.2f}'
    ↪ , end="")
8 if i_episode % 25 == 0:
9     print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps
    ↪ :.2f}')

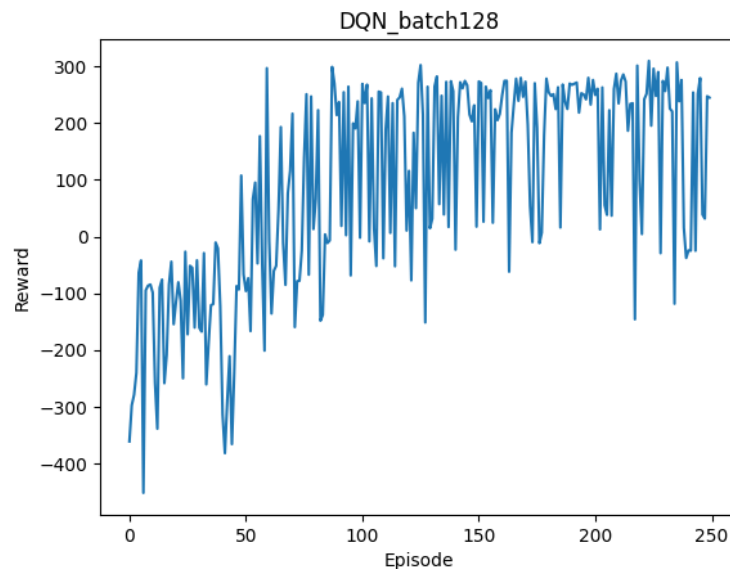
```

Code 14: Epsilon Decay

شکل ۳: نمودار پاداش تجمعی برای  $\text{batch size} = 32$ شکل ۴: نمودار پاداش تجمعی برای  $\text{batch size} = 64$ 

است.

مطابق شکل ۴ و  $\text{batch size} = 64$  در حدود اپیزود ۲۰ یک پاداش منفی بسیار زیاد میگیرد که نشان دهنده یک حادثه خیلی بد برای فضاپیما است. بعد از اپیزود ۵۰ دیده می شود که پاداش ها شروع به صعودی شدن می کنند و البته کمی نوسانات وجود دارد ولی با جلورفتن بهتر وضعیت بهتر می شود. از اپیزود ۱۰۰ پاداش ها حدود ۱۰۰ یا مثبت هستند که نشان می دهد ایجنت در حال پایدارسازی عملکرد و رفتار خود است. می توان گفت در هر سه حالت عملکرد ایجنت با مرور زمان و زیادتر شدن تعداد episode بهبود می یابد و نشان می دهد که الگوریتم درست کار می کند.



شکل ۵: نمودار پاداش تجمعی برای batch size = 128

می توان میانگین هر ۲۵ اپیزود را برای بچ های متفاوت در شکل ۶ تا شکل ۸ مشاهده کرد. مطابق این نمودار ها و به صورت شهودی و با توجه به شکل ۳ تا شکل ۵ می توان گفت که بهترین عملکرد مدل مربوط به حالت batch size = 128 است. این حالت بهترین بالانس را بین دو فاز exploraiton و exploatation ارائه می دهد که در نهایت به بیشترین پاداش ها منجر می شود. در جایگاه ها بعدی به ترتیب batch size = 64 و batch size = 32 قرار دارند که روند بهبود پاداش در آن ها با سرعت کمتری دنبال می شود.

```
Episode 25 Average Reward: -108.03 Epsilon: 0.47
Episode 49 Average Reward: -99.23 Epsilon: 0.22Moviepy - Building video /content/DQN/batch32/eps50/rl-video-episode-0.mp4
Moviepy - Writing video /content/DQN/batch32/eps50/rl-video-episode-0.mp4
/usr/local/lib/python3.10/dist-packages/gym/utils/passive_env_checker.py:233: DeprecationWarning: 'np.bool8' is a deprecated alias for
if not isinstance(terminated, (bool, np.bool8)):
Moviepy - Done !
Moviepy - video ready /content/DQN/batch32/eps50/rl-video-episode-0.mp4
Episode 50 Average Reward: -98.75 Epsilon: 0.22
Episode 75 Average Reward: -83.11 Epsilon: 0.10
Episode 99 Average Reward: -6.59 Epsilon: 0.05Moviepy - Building video /content/DQN/batch32/eps100/rl-video-episode-0.mp4
Moviepy - Writing video /content/DQN/batch32/eps100/rl-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch32/eps100/rl-video-episode-0.mp4
Episode 100 Average Reward: 8.48 Epsilon: 0.05
Episode 125 Average Reward: -61.39 Epsilon: 0.02
Episode 149 Average Reward: 31.99 Epsilon: 0.01Moviepy - Building video /content/DQN/batch32/eps150/rl-video-episode-0.mp4
Moviepy - Writing video /content/DQN/batch32/eps150/rl-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch32/eps150/rl-video-episode-0.mp4
Episode 150 Average Reward: 48.33 Epsilon: 0.01
Episode 175 Average Reward: 142.82 Epsilon: 0.01
Episode 199 Average Reward: 134.62 Epsilon: 0.01Moviepy - Building video /content/DQN/batch32/eps200/rl-video-episode-0.mp4
Moviepy - Writing video /content/DQN/batch32/eps200/rl-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch32/eps200/rl-video-episode-0.mp4
Episode 200 Average Reward: 145.22 Epsilon: 0.01
Episode 225 Average Reward: 214.59 Epsilon: 0.01
Episode 249 Average Reward: 174.53 Epsilon: 0.01Moviepy - Building video /content/DQN/batch32/eps250/rl-video-episode-0.mp4
Moviepy - Writing video /content/DQN/batch32/eps250/rl-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch32/eps250/rl-video-episode-0.mp4
Episode 250 Average Reward: 175.79 Epsilon: 0.01
```

شکل ۶: پاداش تجمعی برای batch size = 32

حال به سراغ ارزیابی این سه حالت به کمک متریک regret می رویم. ابتدا تعریفی از این معیار را با هم مرور کنیم. این معیار برای اندازی گیری اوضاع در شرایط تصمیم گیری و در فضای یادگیری تقویتی به کار می رود. این معیار



```
Episode 25 Average Reward: -157.51 Epsilon: 0.47
Episode 49 Average Reward: -178.48 Epsilon: 0.22Moviepy - Building video /content/DQN/batch64/eps50/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch64/eps50/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch64/eps50/r1-video-episode-0.mp4
Episode 50 Average Reward: -175.24 Epsilon: 0.22
Episode 75 Average Reward: 6.67 Epsilon: 0.10
Episode 99 Average Reward: 37.77 Epsilon: 0.05Moviepy - Building video /content/DQN/batch64/eps100/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch64/eps100/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch64/eps100/r1-video-episode-0.mp4
Episode 100 Average Reward: 49.86 Epsilon: 0.05
Episode 125 Average Reward: 113.18 Epsilon: 0.02
Episode 149 Average Reward: 170.41 Epsilon: 0.01Moviepy - Building video /content/DQN/batch64/eps150/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch64/eps150/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch64/eps150/r1-video-episode-0.mp4
Episode 150 Average Reward: 162.35 Epsilon: 0.01
Episode 175 Average Reward: 92.24 Epsilon: 0.01
Episode 199 Average Reward: 169.26 Epsilon: 0.01Moviepy - Building video /content/DQN/batch64/eps200/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch64/eps200/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch64/eps200/r1-video-episode-0.mp4
Episode 200 Average Reward: 177.33 Epsilon: 0.01
Episode 225 Average Reward: 118.35 Epsilon: 0.01
Episode 249 Average Reward: 75.09 Epsilon: 0.01Moviepy - Building video /content/DQN/batch64/eps250/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch64/eps250/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch64/eps250/r1-video-episode-0.mp4
Episode 250 Average Reward: 76.11 Epsilon: 0.01
```

شکل ۷: پاداش تجمعی برای batch size = 64

```
Episode 25 Average Reward: -168.00 Epsilon: 0.47
Episode 49 Average Reward: -143.30 Epsilon: 0.22Moviepy - Building video /content/DQN/batch128/eps50/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch128/eps50/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch128/eps50/r1-video-episode-0.mp4
Episode 50 Average Reward: -144.90 Epsilon: 0.22
Episode 75 Average Reward: -1.75 Epsilon: 0.10
Episode 99 Average Reward: 106.53 Epsilon: 0.05Moviepy - Building video /content/DQN/batch128/eps100/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch128/eps100/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch128/eps100/r1-video-episode-0.mp4
Episode 100 Average Reward: 107.44 Epsilon: 0.05
Episode 125 Average Reward: 142.99 Epsilon: 0.02
Episode 149 Average Reward: 182.01 Epsilon: 0.01Moviepy - Building video /content/DQN/batch128/eps150/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch128/eps150/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch128/eps150/r1-video-episode-0.mp4
Episode 150 Average Reward: 172.46 Epsilon: 0.01
Episode 175 Average Reward: 199.31 Epsilon: 0.01
Episode 199 Average Reward: 218.05 Epsilon: 0.01Moviepy - Building video /content/DQN/batch128/eps200/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch128/eps200/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch128/eps200/r1-video-episode-0.mp4
Episode 200 Average Reward: 218.28 Epsilon: 0.01
Episode 225 Average Reward: 186.54 Epsilon: 0.01
Episode 249 Average Reward: 159.23 Epsilon: 0.01Moviepy - Building video /content/DQN/batch128/eps250/r1-video-episode-0.mp4.
Moviepy - Writing video /content/DQN/batch128/eps250/r1-video-episode-0.mp4
Moviepy - Done !
Moviepy - video ready /content/DQN/batch128/eps250/r1-video-episode-0.mp4
Episode 250 Average Reward: 161.21 Epsilon: 0.01
```

شکل ۸: پاداش تجمعی برای batch size = 128

تفاوت بین پاداشی که توسط ایجنت دریافت شده و پاداشی که دریافت می شد، اگر ایجنت بهترین تصمیم را در هر اپیزود می گرفت ارزیابی می کند. پس می توان آن را به صورت رابطه ۱ فرموله کرد.

$$R_t = \mu^* - \mu_a \quad (1)$$

در رابطه ۱  $\mu^*$  پاداشی است که اگر ایجنت بهترین اکشن را انتخاب می کرد دریافت می کرد. در طرف دیگر  $\mu_a$  پاداشی است که ایجنت با انتخاب کردن اکشن  $a$  در زمان  $t$  آن را دریافت کرده است.

حال اگر ما رابطه ۱ را در تمامی اپیزود ها حساب کنیم می توانیم regret تجمعی برای ایجنت را بدست آوریم.

در واقع می توان گفت regret هزینه انتخاب نکردن بهترین تصمیم در هر اپیزود را به ما نشان می دهد.

مطابق نتایج حالت batch size = 32 در ابتدا regret متوسط است، سپس این معیار به صورت واضحی کم میشود چون reward ها در حال افزایش هستند (حدود اپیزود ۵۰) و در نهایت در اپیزود های بالاتر این معیار به کمترین حالت خود میرسد از آنجایی که پاداش ها در حال افزایش هستند.

برای حالت batch size = 128 ما کمترین میزان regret را در فاز های مختلف شاهد هستیم چون بیشتری پاداش هارا برای این حالت دریافت می کنیم و می توان اینطور توصیف کرد که در این حالت ایجنت به بهترین عملکرد خود نزدیک است. پس با توجه به این معیار نیز بهترین گزینه ما حالت batch size = 128 است.

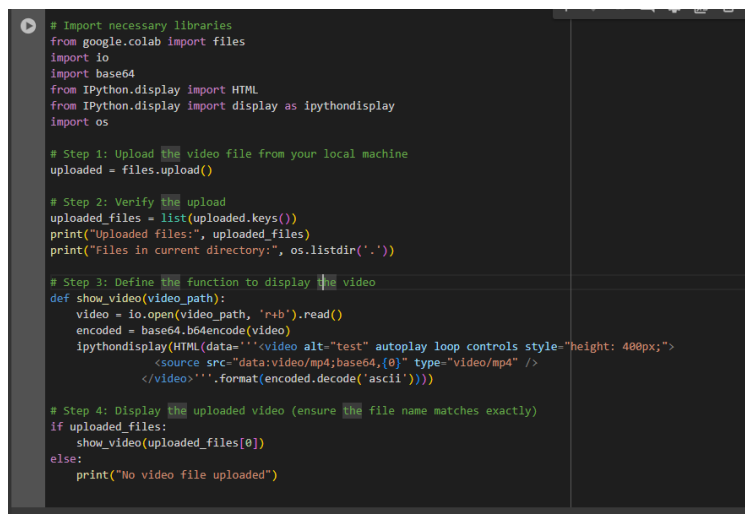




در نهایت این بخش به سراغ تهیه فیلم ها می رویم که در روند آموزش آن هارا ذخیره کرده ایم. توجه شود که تمامی فیلم های مربوط به هر سه بچ در فایل ارائه آمده است. برای حالت  $batch\ size = 128$  این ویدیو ها روی فایل google collab به صورت **شکل ۹** قرار گرفته اند.



شکل ۹: ویدیویی از ایجنت در محیط به ازای  $batch\ size = 128$  در اپیزود ۲۵۰



شکل ۱۰: ساختن ویدیو از ایجنت در محیط  $batch\ size = 128$  در اپیزود ۲۵۰

این ویدیو ها به ازای تمامی اپیزود ها در فایل google collab قرار گرفته اند.

## ۳.۲ قسمت سوم

باید در این قسمت مدل DDQN را به جای مدل DQN جاگذاری کنیم. می توان گفت تفاوت اصلی این دو مدل در نحوه بروز کرد Q value ها و شبکه هدف می باشد. به طور دقیق تر DQN از شبکه استفاده می کند تا Q value ها را تخمین بزند اما DDQN از شبکه موجود اکشن را انتخاب می کند و Q value های مرحله بعد را به کمک شبکه هدف تخمین می زند. این تفاوت در برنامه ۱۵ و برنامه ۱۶ قابل مشاهده است.



```
1 q_targets_next = self.target_net(next_state_batch).detach().max(1)[0].unsqueeze(1)
2 q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch))
```

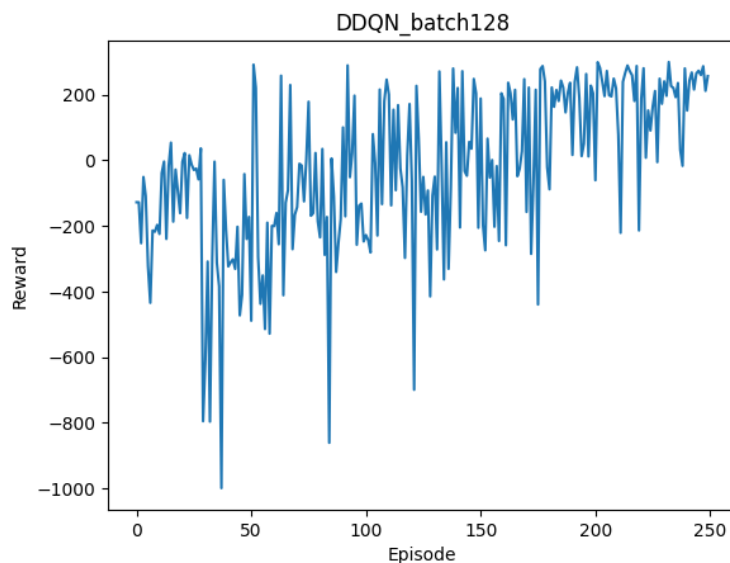
Code 15: DQN Agent

```
1 next_action_batch = torch.argmax(self.value_net(next_state_batch), dim=1, keepdim=True)
2 q_targets_next = self.target_net(next_state_batch).gather(1, next_action_batch).detach()
3 q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch))
```

Code 16: DDQN Agent

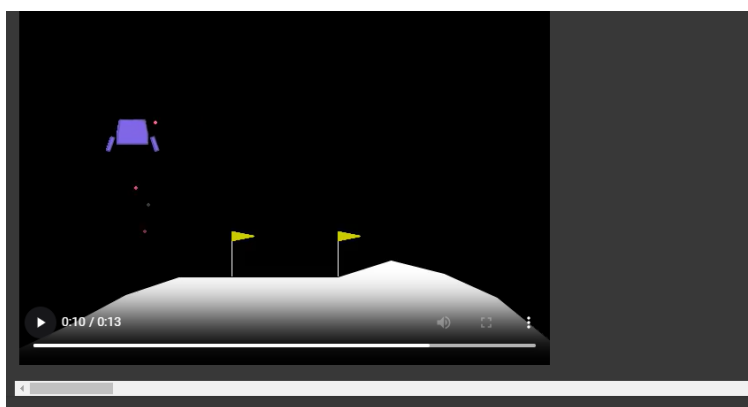
مشاهده می شود که DDQN یک next action batch دارد که انتخاب اکشن با کمک شبکه فعلی را انجام می دهد که در DQN وجود ندارد. خوبی این روش این است که از بایاس های احتمالی در شبکه DQN جلوگیری می کند چون گام های انتخاب اکشن و تخمین Q value با هم تفاوت دارد. در نهایت این روش می تواند به پاسخ هایی پایدارتر و با اتکاتری را ارائه دهد.

نتایج شبیه سازی با کمک DDQN در حالت batch size = 128 در شکل ۱۱ آمده است. در این حالت دیده می شود



شکل ۱۱: نمودار پاداش تجمعی به کمک DDQN در حالت batch size = 128

که پاداش ها بسیار پایدار تر از حالت قبل هستند، یعنی پایداری ایجنت در دریافت پاداش های بیشتر بهتر شده است. برای تهیه ویدیو مانند بخش قبل از دستورات مربوط به نمایش ویدیو در محیط google collab استفاده می کنیم و به صورت گفته شده برای دو حالت اپیزود ۱۰۰ و ۲۵۰ ویدیو را نشان می دهیم. قابل ذکر است این ویدیو ها در فایل ارسالی قرار دارند.



شکل ۱۲: ویدیویی از ایجنت در محیط به ازای  $\text{batch size} = 128$  در اپیزود ۲۵۰