

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

یادگیری ماشین پاسخ مینی پروژه دوم

نام و نام خانوادگی: سید محمد حسینی
شماره دانشجویی: ۹۹۰۱۳۹۹
تاریخ: مهرماه ۱۴۰۲

GitHub

Drive Google



فهرست مطالب

۵	۱ پرسش ۲
۵	۱.۱ قسمت اول
۶	۲.۱ قسمت دوم
۹	۱.۲.۱ تعریف متغیرها و مقادیر اولیه
۹	۲.۲.۱ متد <code>take_action</code>
۹	۳.۲.۱ متد <code>update_params</code>
۱۱	۴.۲.۱ Training
۱۱	۵.۲.۱ نتایج
۱۳	۳.۱ قسمت سوم
۱۳	۱.۳.۱ DQN
۱۴	۲.۳.۱ DDQN



فهرست تصاویر

۱۲	نتایج به ازای تمام اپیزودها	۱
۱۳	نتایج به ازای هر ۵۰ اپیزود	۲
۱۵	مقایسه DQN و DDQN	۳
۱۶	مقایسه DQN و DDQN	۴



فهرست جداول



فهرست برنامه‌ها

۶ Network Neural Deep	۱
۷ Class Agent	۲
۹ Train	۳



۱ پرسش ۲

۱.۱ قسمت اول

lunar lander فضایی است که برای فرود آمدن بر سطح ماه طراحی شده است و مأموریت آن فرود ایمن از مدار ماه به سطح ماه است. هدف اصلی این فضاپیما شامل طی کردن مسیر دقیق فرود و فرود ایمن و نرم است که نیازمند سیستم‌های راهنمایی پیشرفته، ناوبری دقیق و طراحی ساختاری قوی است. سیستم پیشران که شامل موتورهای اصلی فرود و موتورهای کنترل وضعیت جانبی و کوچکتر است که نقش مهمی در کنترل سرعت فرود و جهت‌گیری ایفا می‌کند. lunar lander همچنین دارای پایه‌های فرود برای کاهش اثر ضربه و تضمین پایداری هستند.

یادگیری تقویتی (RL) می‌تواند برای مسئله کنترل یک lunar lander استفاده شود. در این زمینه، RL شامل آموزش یک Agent برای یادگیری یک Policy بهینه برای فرود ایمن و کارآمد م lunar lander است. اجزای اصلی یک مسئله RL شامل فضای حالت، فضای عمل و سیستم پاداش است که این موارد به شرح زیر می‌باشند.

فضای حالت

فضای حالت تمام حالات ممکن را که فرودگر قمری می‌تواند در آن‌ها باشد، نشان می‌دهد. برای یک lunar lander فضای حالت می‌تواند شامل موارد زیر باشد:

- مختصات عمودی و افقی فرودگر.
 - مولفه‌های سرعت در جهت‌های افقی و عمودی.
 - زاویه فرودگر نسبت به عمود.
 - سرعت زاویه‌ای فرودگر.
 - وضعیت پایه‌های فرودگر
- هر حالت نمایانگر شرایط فرودگر در یک گام زمانی معین است.

فضای عمل

فضای عمل تمام Actions ممکن را که Agent می‌تواند انجام دهد، نشان می‌دهد. برای یک lunar lander فضای عمل می‌تواند شامل موارد زیر باشد:

- توان موتور اصلی.
- توان موتور چپ.
- توان موتور راست.
- انجام ندادن هیچ کاری.

این Actions نیروی پیشران و جهت‌گیری فرودگر را کنترل می‌کنند.

سیستم پاداش



سیستم پاداش بازخوردی به Agent بر اساس Actions آن و تغییرات حالت ناشی از آن‌ها ارائه می‌دهد. برای یک lunar lander سیستم پاداش ممکن است به این صورت طراحی شود:

- پاداش مثبت برای دستیابی به فرود ایمن با سرعت کم.
- پاداش منفی برای فرودهای سخت یا سقوط‌ها.
- پاداش منفی برای مصرف سوخت زیاد.
- پاداش مثبت برای حفظ جهت‌گیری پایدار.
- پاداش منفی کوچک برای هر گام زمانی برای تشویق فرود سریع‌تر.

سیستم پاداش Agent را تشویق می‌کند تا هایی Policy را بیاموزد که منجر به فرودهای ایمن، کارآمد و پایدار شود. [۱]

۲.۱ قسمت دوم

ابتدا به بررسی بلاک‌های کد می‌پردازیم.

```
۱ # DQN
۲ import torch.nn as nn
۳ import torch.nn.functional as F
۴
۵ class DeepQNetwork(nn.Module):
۶     def __init__(self, state_size, action_size) -> None:
۷         super(DeepQNetwork, self).__init__()
۸         # TODO: define the architecture
۹         # NOTE: input=observation/state, output=action
۱0        net_list = nn.ModuleList([
۱1            torch.nn.Linear(state_size, 512),
۱2            torch.nn.ReLU(),
۱3            torch.nn.LayerNorm(512),
۱4            torch.nn.Dropout(0.1),
۱5            torch.nn.Linear(512, 512),
۱6            torch.nn.ReLU(),
۱7            torch.nn.LayerNorm(512),
۱8            torch.nn.Dropout(0.1),
۱9            torch.nn.Linear(512, 512),
۲0            torch.nn.ReLU(),
۲1            torch.nn.Linear(512, action_size)
۲2        ])
۲3        self.net = torch.nn.Sequential(*net_list).to(device)
۲۴
۲۵     def forward(self, x):
۲۶         # TODO: forward propagation
```



```
27 # NOTE: use ReLu for activation function in all layers
28 # NOTE: last layer has no activation function (predict action)
29 # ReLU is created in init, no need here
30 x.to(device)
31 x = self.net(x)
32 return x
33
```

Code : ۱ Deep Neural Network

کد فوق یک شبکه عصبی عمیق را ایجاد می‌کند که به موجب آن ۴ لایه خطی به همراه ۳ لایه غیر خطی ایجاد شده است تا ماتریس پاداش را پیشبینی کند.

```
1 # DQN agent
2 import numpy as np
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7
8 class DQNAgent():
9     # NOTE: DON'T change initial values
10     def __init__(self, state_size, action_size, batch_size,
11                 gamma=0.99, buffer_size=25000, alpha=1e-4):
12         # network parameter
13         self.state_size = state_size
14         self.action_size = action_size
15
16         # hyperparameters
17         self.batch_size = batch_size
18         self.gamma = gamma
19
20         # experience replay
21         self.experience_replay = ExperienceReplay(buffer_size)
22
23         # network
24         self.value_net = DeepQNetwork(state_size, action_size).to(device)
25
26         # optimizer
27         # TODO: create adam for optimizing network's parameter (learning rate=alpha)
28         self.optimizer = optim.Adam(self.value_net.parameters(), lr=alpha)
29
30     def take_action(self, state, eps=0.0):
31         # TODO: take action using e-greedy policy
32         # NOTE: takes action using the greedy policy with a probability of -1 and a random action
33         # NOTE: with a probability of
```




```
۳۴     self.value_net.eval()
۳۵     if len(state) != 8 :
۳۶         state = state[0]
۳۷     rand_eps = random.random()
۳۸     if rand_eps > eps:
۳۹         with torch.no_grad():
۴۰             # print(state)
۴۱             return torch.argmax(self.value_net(torch.tensor(state).to(device))).detach().cpu
().numpy()
۴۲     else:
۴۳         return np.random.randint(0, self.action_size)
۴۴
۴۵ def update_params(self):
۴۶     if len(self.experience_replay) < self.batch_size:
۴۷         return
۴۸     # transition batch
۴۹     batch = Transition(*zip(*self.experience_replay.sample(self.batch_size)))
۵۰
۵۱     temp = []
۵۲     for indx in range(len(batch.state)):
۵۳         if len(batch.state[indx]) != 8:
۵۴             temp.append(batch.state[indx][0])
۵۵         else:
۵۶             temp.append(batch.state[indx])
۵۷     state_batch = torch.from_numpy(np.vstack(temp)).float().to(device) # [8, 8]
۵۸     action_batch = torch.tensor(np.vstack(batch.action)).long().to(device) # [8, 1]
۵۹     next_state_batch = torch.from_numpy(np.vstack(batch.next_state)).float().to(device) # [8,
8]
۶۰     reward_batch = torch.tensor(np.vstack(batch.reward)).float().to(device) # [8, 1]
۶۱     done_batch = torch.tensor(np.vstack(batch.done)).to(device)
۶۲
۶۳     # calculate loss w.r.t DQN algorithm
۶۴     self.value_net.train()
۶۵     # STEP1
۶۶     q_expected = self.value_net(state_batch).gather(1, action_batch)
۶۷     # TODO: compute the expected Q values [y]
۶۸     # STEP2
۶۹     # TODO: compute Q values [Q(s_t, a)]
۷۰     q_targets_next = self.value_net(next_state_batch).detach().max(1)[0].unsqueeze(1)
۷۱     q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch*1))
۷۲     # STEP3
۷۳     # TODO: compute mse loss
۷۴     loss = nn.functional.mse_loss(q_expected, q_targets)
۷۵     # TODO: optimize the model
۷۶     # NOTE: DON'T forget to set the gradients to zeros
```



```

۷۷     self.optimizer.zero_grad()
۷۸     loss.backward()
۷۹     self.optimizer.step()
۸۰
۸۱     def save(self, fname):
۸۲         # TODO: save checkpoint
۸۳         torch.save(self.value_net, fname)
۸۴
۸۵     def load(self, fname, device):
۸۶         # TODO: load checkpoint
۸۷         self.value_net = torch.load(fname, device)
۸۸

```

Code : ۲ Agent Class

در این بخش، کلاس DQNAgent تعریف می‌شود که Agent مربوط به DQN را پیاده‌سازی می‌کند. این کلاس مسئول مدیریت شبکه عصبی، اجرای Policy ϵ -Greedy، به‌روزرسانی پارامترهای شبکه و ذخیره و بارگذاری مدل است.

۱.۲.۱ تعریف متغیرها و مقادیر اولیه

در ابتدای کلاس، مقادیر اولیه و پارامترهای شبکه تعریف می‌شوند. این شامل اندازه State ($state_size$)، اندازه Action ($action_size$)، اندازه Batch ($batch_size$)، ($gamma$)، اندازه Buffer ($buffer_size$) و نرخ یادگیری ($alpha$) است. در ادامه شبکه عصبی (DeepQNetwork) برای پیش‌بینی مقادیر Q و بهینه‌ساز Adam با نرخ یادگیری α تعریف می‌شوند.

۲.۲.۱ متد take_action

این تابع برای انتخاب اقدام با استفاده از Policy ϵ -Greedy استفاده می‌شود. با احتمال $1 - \epsilon$ ، اقدام بهینه انتخاب می‌شود و با احتمال ϵ ، یک اقدام تصادفی انتخاب می‌شود.

۳.۲.۱ متد update_params

این تابع پارامترهای شبکه عصبی را به‌روزرسانی می‌کند. اگر تعداد تجربیات کمتر از $batch_size$ باشد، تابع متوقف می‌شود. در ادامه تجربیات نمونه‌گیری شده و به فرمت مناسب برای استفاده در شبکه عصبی تبدیل می‌شوند.

حال مقادیر Q مورد انتظار ($q_expected$) و سپس مقادیر هدف Q ($q_targets$) محاسبه می‌شوند. سپس با استفاده از (MSE) تابع هزینه محاسبه شده و مدل بهینه‌سازی می‌شود.

```

۱ # training phase
۲
۳ # TODO: create agent

```



```
۴ agent = DQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
۵
۶ crs = np.zeros(n_episodes) # cummulative rewards
۷ crs_recent = deque(maxlen=25) # recent cummulative rewards
۸
۹ # training loop
۱۰ for i_episode in range(1, n_episodes+1):
۱۱     # TODO: initialize the environment and state
۱۲     if i_episode % 50 == 0:
۱۳         env = RecordVideo(gym.make("LunarLander-v2"), f"./DQN/batch{BATCH_SIZE}/eps{i_episode}")
۱۴     else:
۱۵         env = gym.make("LunarLander-v2")
۱۶     state = env.reset()
۱۷     done = False
۱۸     cr = 0 # episode cummulative rewards
۱۹     while not done:
۲۰         env.render()
۲۱         # TODO: select and perform an action
۲۲         action = agent.take_action(state, eps)
۲۳     #         print(action)
۲۴         # Modify the unpacking to handle the extra value if present
۲۵         result = env.step(action)
۲۶         if len(result) == 5:
۲۷             next_state, reward, done, truncated, info = result
۲۸         else:
۲۹             next_state, reward, done, info = result
۳۰         # TODO: store transition in experience replay
۳۱         agent.experience_replay.store_trans(state, action, next_state, reward, done)
۳۲         # TODO: update agent
۳۳         agent.update_params()
۳۴         # TODO: update current state and episode cummulative rewards
۳۵         # print("next" ,next_state)
۳۶         state = next_state
۳۷         cr += reward
۳۸
۳۹     # TODO: decay epsilon
۴۰     eps = eps * eps_decay_rate
۴۱     eps = max(eps, eps_end)
۴۲     # TODO: update current cummulative rewards and recent cummulative rewards
۴۳     crs[i_episode - 1] = cr
۴۴     crs_recent.append(cr)
۴۵     # TODO: save agent every 50 episodes
۴۶     if i_episode % 50 == 0:
۴۷         agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
۴۸
```



```

۴۹ # print logs
۵۰ print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(i_episode, np.mean(
    crs_recent), eps), end="")
۵۱ if i_episode % 25 == 0:
۵۲     print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(i_episode, np.mean(
        crs_recent), eps))
۵۳

```

Code :۳ Train

۴.۲.۱ Training

در این بخش، کدهای مربوط به بخش Training توضیح داده می‌شود. این فاز شامل ایجاد Agent، اجرای حلقه آموزش، انتخاب Actions، به‌روزرسانی پارامترها، و ذخیره مدل است.

در گام نخست Agent مد نظر ما با مقادیر اولیه متناسب با مسئله ایجاد می‌شود. در ادامه متغیر crs برای ذخیره پاداش‌های تجمعی هر اپیزود و crs_recent برای ذخیره پاداش‌های تجمعی اخیر تعریف می‌شود.

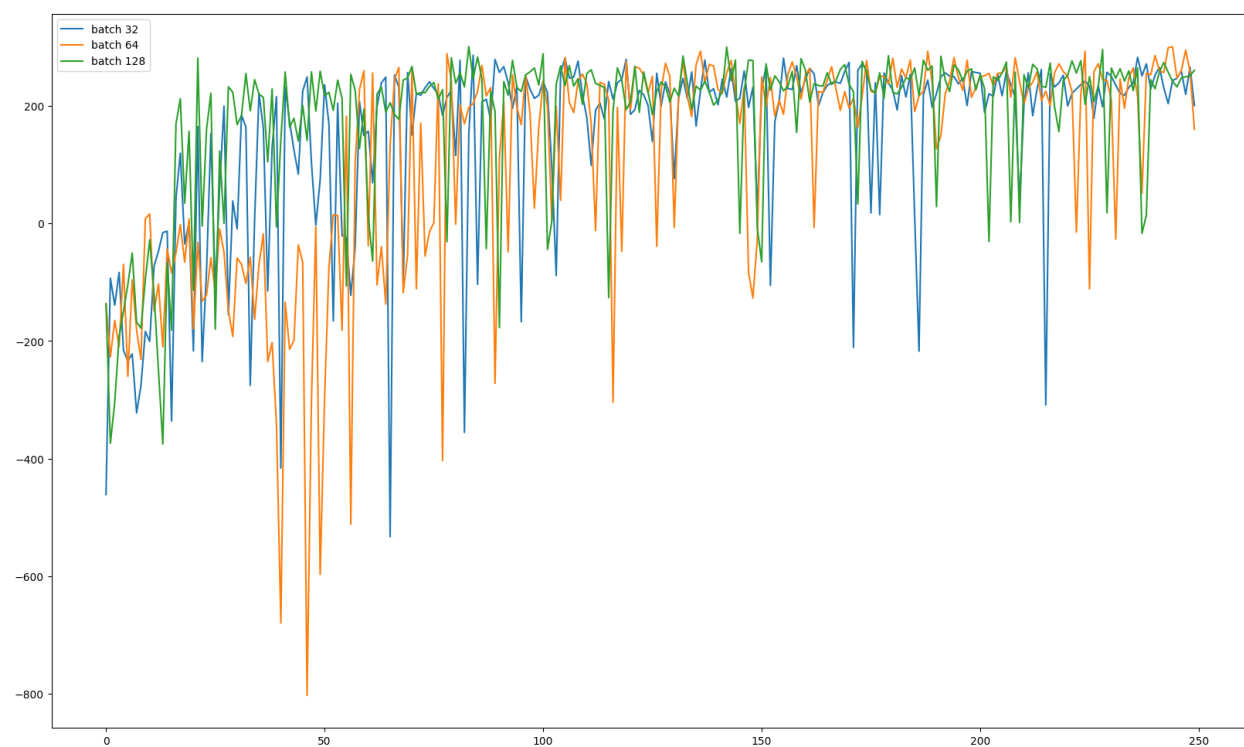
حلقه آموزش برای تعداد اپیزودهای مشخص شده (n_episodes) اجرا می‌شود تا Agent پاسخ بهینه به مسئله را پیدا کند. در این حلقه آموزش به ازای هر اپیزود، محیط و حالت اولیه مقداردهی می‌شوند. هر ۵۰ اپیزود، محیط برای RecordVideo تنظیم می‌شود.

تا زمانی که اپیزود به پایان نرسیده است، حلقه steps اجرا می‌شود. در هر گام، محیط رندر می‌شود و Agent یک اقدام انتخاب می‌کند. نتیجه اقدام انجام شده دریافت می‌شود و انتقال در تجربه‌های Agent ذخیره می‌شود. پارامترهای Agent به‌روزرسانی می‌شوند و حالت فعلی و پاداش‌های تجمعی اپیزود به‌روز می‌شوند.

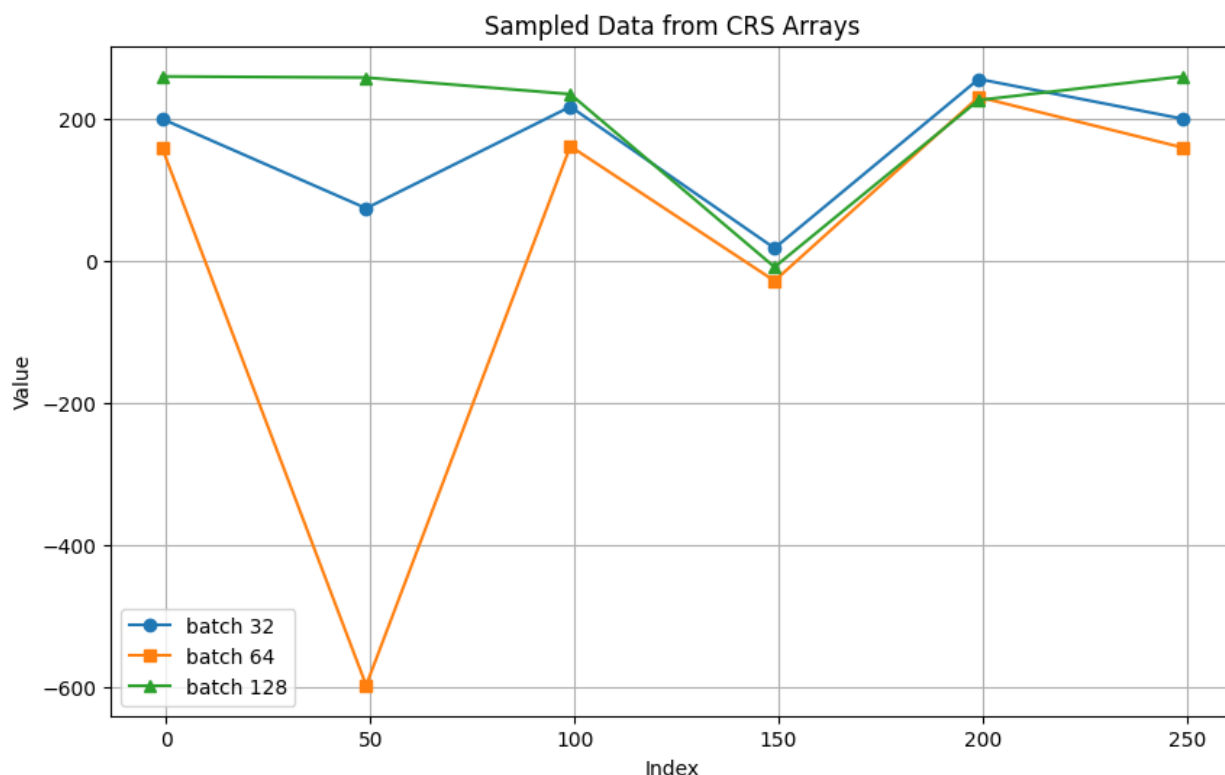
پس از پایان اپیزود، مقدار epsilon کاهش می‌یابد تا به تدریج Agent کمتر به Actions تصادفی روی بیاورد. پاداش‌های تجمعی اپیزود جاری و پاداش‌های تجمعی اخیر به‌روز می‌شوند. هر ۵۰ اپیزود، مدل Agent ذخیره می‌شود.

۵.۲.۱ نتایج

نمودارهای زیر برای نمایش نتایج حاصله از آموزش مدل رسم شده‌اند.



شکل ۱: نتایج به ازای تمام اپیزودها



شکل ۲: نتایج به ازای هر ۵۰ اپیزود

همانطور که انتظار می‌رفت نتایج به ازای بچ ۱۲۸ زودتر همگرا شده است اما به ازای بچ ۶۴ نتایج همگرایی ضعیفتری داشته است نسبت به حالتی که بچ ۳۲ بوده است. بنابراین بهترین حالت مربوط به بچ ۱۲۸ بوده و به منظور توجه به سرعت همگرایی در اپیزود ۵۰ بهترین مدل انتخاب می‌شود. کد اجرایی به منظور ذخیره ویدیو از آموزش Agent همراه با یک ارور بوده که امکان حل آن وجود نداشت. ارور در فایل نوتبوک موجود می‌باشد.

۳.۱ قسمت سوم

در یادگیری تقویتی، DQN و DDQN دو روش محبوب برای یادگیری مقادیر-Q هستند. هر دو روش به منظور بهبود عملکرد Agent در انتخاب بهترین Actions در یک محیط مشخص توسعه داده شده‌اند. در اینجا تفاوت‌ها و شباهت‌های کلیدی بین DQN و DDQN را بررسی می‌کنیم.

۱.۳.۱ DQN

- **معماری:** DQN از یک شبکه عصبی عمیق برای تقریب مقادیر-Q استفاده می‌کند. این شبکه عصبی ورودی‌هایی شامل حالت‌های محیط را دریافت کرده و مقادیر-Q مربوط به هر اقدام ممکن را تولید می‌کند.

- به روزرسانی: مقادیر-Q در DQN، از یک شبکه هدف استفاده می شود که کپی ای از شبکه اصلی (Q-Network) است و در فواصل زمانی منظم به روزرسانی می شود. این کار به پایداری یادگیری کمک می کند.

- معادله به روزرسانی: معادله به روزرسانی مقادیر-Q در DQN به صورت زیر است:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left(r + \gamma \max_{a'} Q'(s', a') - Q(s, a) \right)$$

که در آن $Q'(s', a')$ مقدار Q پیش بینی شده توسط شبکه هدف است.

- Over fitting در DQN، به دلیل استفاده از $\max_{a'} Q(s', a')$ برای به روزرسانی، مقادیر-Q احتمال Over fitting وجود دارد که می تواند منجر به پایداری کمتر در یادگیری شود.

۲.۳.۱ DDQN

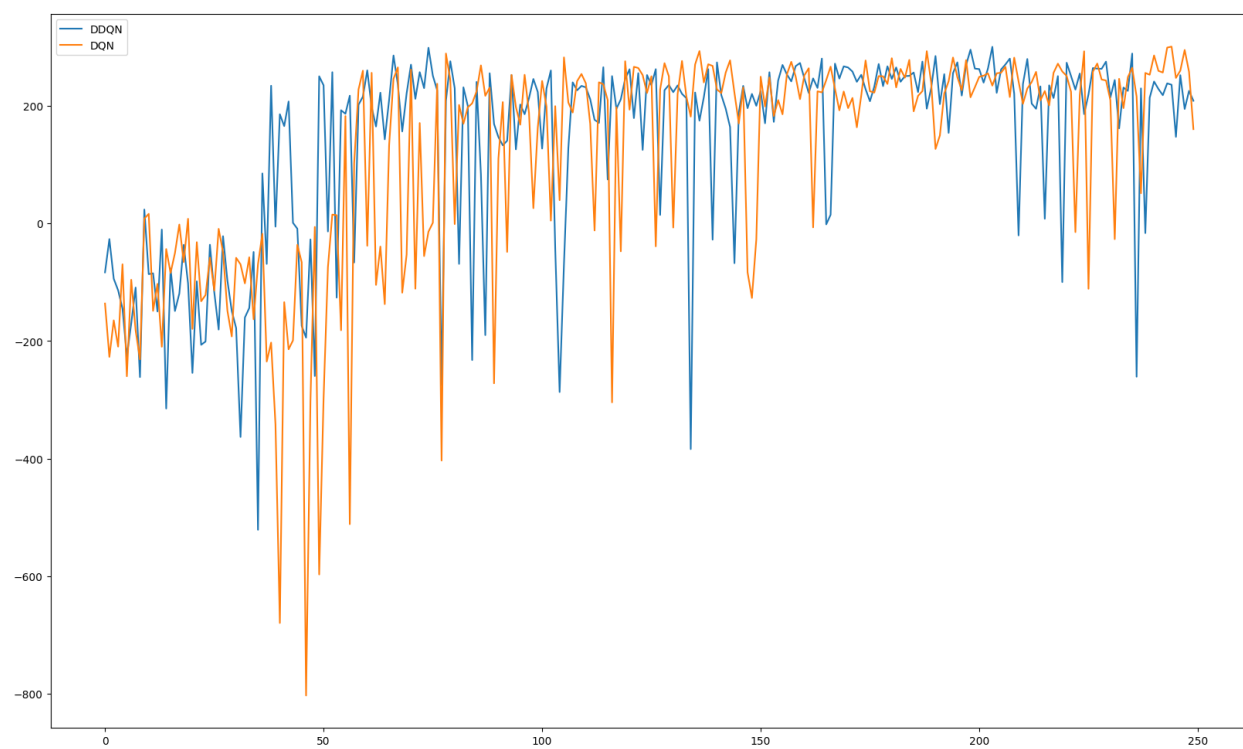
- معماری: معماری شبکه عصبی در DDQN مشابه DQN است، با این تفاوت که در DDQN دو شبکه Q به صورت جداگانه به روزرسانی می شوند: یک شبکه Q اصلی و یک شبکه Q هدف.

- به روزرسانی: مقادیر-Q در DDQN، به روزرسانی مقادیر-Q به گونه ای اصلاح شده است که از Over fitting جلوگیری شود. این کار با استفاده از شبکه Q اصلی برای انتخاب Actions و شبکه Q هدف برای محاسبه مقادیر-Q انجام می شود.

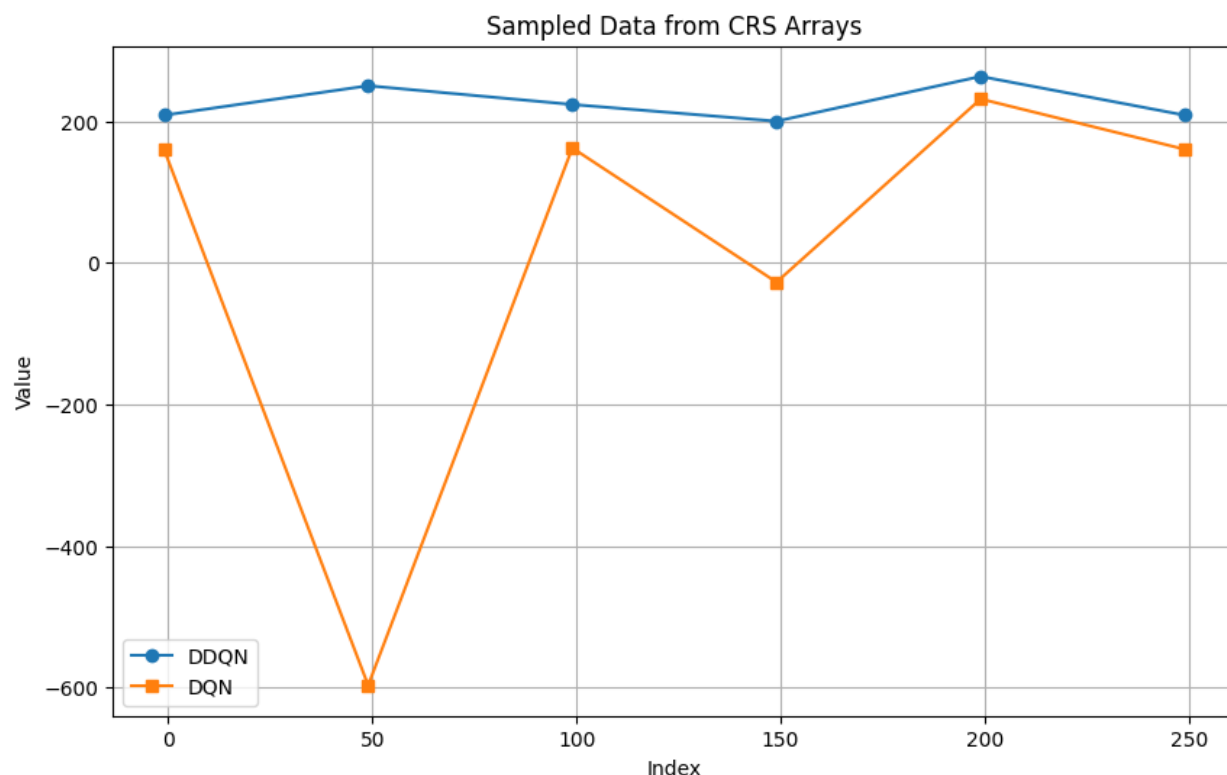
- معادله به روزرسانی: معادله به روزرسانی مقادیر-Q در DDQN به صورت زیر است:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left(r + \gamma Q'(s', \arg \max_{a'} Q(s', a')) - Q(s, a) \right)$$

که در آن $Q(s', a')$ مقدار Q پیش بینی شده توسط شبکه اصلی است و $Q'(s', \arg \max_{a'} Q(s', a'))$ مقدار Q پیش بینی شده توسط شبکه هدف با استفاده از اقدام انتخاب شده توسط شبکه Q اصلی است.



شکل ۳: مقایسه DDQN و DQN



شکل ۴: مقایسه DDQN و DQN

همانطور که انتظار می‌رفت مدل DDQN نتیجه بهتر و پایداری کسب کرده است که این مسئله ناشی از تغییر روند آموزش شبکه عصبی می‌باشد. پایداری مدل DDQN در فرایند همگرایی و نتیجه کلی Agent تاثیر مهمی داشته است.

مراجع

[۱] Xu. Sophia course, science data for project Final Accessed: ۲۰۲۴-۰۷-۰۵.