

به نام خدا



## مینی پروژه چهارم یادگیری ماشین

استاد: دکتر مهدی علیاری

محمدقاسم امامی مقدم

شماره دانشجویی: ۴۰۲۰۲۳۱۴

پوشه گیت هاب مینی پروژه - پوشه گوگل کولب مینی پروژه

اردیبهشت ۱۴۰۳

## فهرست مطالب

۱	..... سوال اول	۱
۱	..... بخش اول	۱.۱
۲	..... بخش دوم	۲.۱
۵	..... بخش سوم	۳.۱

## فهرست تصاویر

۱	..... صورت کلی بازی موجود در اینترنت	۱
۲	..... مقایسه دو روش	۲
۵	..... DDQN	۳
۵	..... دسته های ۳۲ عددی	۴
۶	..... دسته های ۶۴ عددی	۵
۶	..... دسته های ۱۲۸ عددی	۶

# ۱ سوال اول

## ۱.۱ بخش اول

هدف بازی Lunar Lander هدایت یک سفینه میباشد که طوری که باید در محدوده مورد نظر بنشیند. در واقع باید از محل شروع بازی و رهاسازی سفینه با مولفه های موجود و در دسترس سفینه را در حالت بهینه کنترل کرد و در مقصد نشانند.



شکل ۱: صورت کلی بازی موجود در اینترنت

متغیر های حالت مربوط به سفینه

۱. موقعیت عمودی  $x$

۲. موقعیت افقی  $y$

۳. سرعت عمودی  $v_y$

۴. سرعت افقی  $v_x$

۵. زاویه  $\theta$

۶. سرعت زاویه ای  $w$

۷. مقادیر بولین ←

□ تماس پایه چپ با سطح left leg contact

□ تماس پایه راست با سطح right leg contact

دقت شود که مقادیر مربوط به موقعیت های سفینه به صورت نسبی در نظر گرفته میشوند بدین صورت که محل فرود را نقطه (۰ و ۰) در نظر میگیریم و موقعیت عمودی و افقی را نسبت به آن میسنجیم. در این سیستم نیز ورودی های مربوط به سفینه قابل کنترل توسط ما هستند. در واقع با این مولفه ها میتوان تمامی متغیر های حالت مربوط به مسله را کنترل نمود. در هر لحظه موقعیت سفینه برای ما مدنظر می باشد. فضای اکشن مربوط به سفینه:

۱. هیچ کار انجام ندادن

۲. روشن کردن موتور اصلی

۳. روشن کردن موتور سمت راستی

۴. روشن کردن موتور سمت چپی

برای این سوال می بایست یک سیستم امتیاز دهی در نظر گرفته شود تا بتوان عملکرد سیستم را بهبود بخشید و هر مرحله نسبت به مرحله قبلی عملکرد بهتری داشته باشد. در ابتدا برای عملکرد های مثبت و مثر ثمر امتیازات بالا و مثبت در نظر گرفته میشود.

۱. ساکن شدن سفینه : ۱۰۰ امتیاز مثبت

۲. تماس هر پایه با زمین : ۱۰ امتیاز مثبت

۳. حرکت به سمت محدوده مطلوب : ۱۰۰ تا ۱۴۰ امتیاز مثبت

۴. حل سوال و جاگیری در محدوده مطلوب : ۲۰۰ امتیاز مثبت

در طرف مقابل امتیازات منفی نیز برای سیستم در نظر گرفته میشود:

۱. روشن کردن موتور اصلی(بالابرنده) : ۳۰۰ امتیاز منفی

۲. روشن کردن موتور های جانبی : ۰۳۰۰ امتیاز منفی

۳. برخورد با زمین در اثر سقوط یا در جهت اشتباه (عدم قود روی پایه های سفینه) : ۱۰۰ امتیاز منفی

حال با توجه به شرایط داده شده و امتیازات اختصاص داده شده به سیستم در ابتدا سفینه از یک نقطه اولیه قرارند اغار می شود. در ابتدا شرایط اولیه وارده به سیستم به صورت رندوم در نظر گرفته میشود. تنها نیروی وارده به سیستم نیروی جاذبه ی وارده به سیستم میباشد . فرایند امتیاز دهی و سوال در چند حالت به اتمام میرسد.

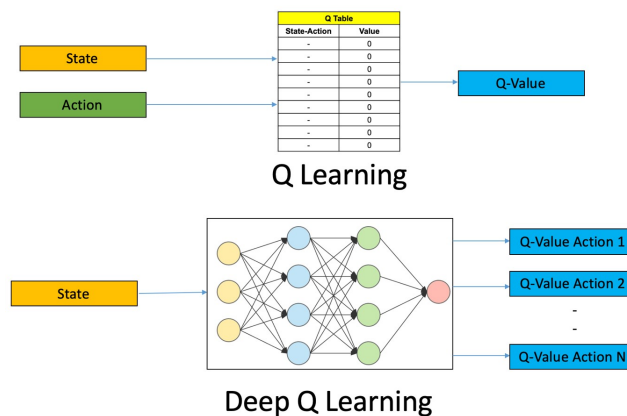
۱. سقوط کردن سفینه و برخورد با سطح

۲. خارج شدن سفینه از محدوده سوال

۳. رسیدن به حالت نهایی

## ۲.۱ بخش دوم

دقت شود در مورد سوالاتی که فضای سوال مربوطه کوچک باشد و خیلی بزرگ نباشد از الگوریتم Q-Learning استفاده میکنیم و همان طور که در کلاس مورد بررسی قرار گرفت ماتریس مربوط به Q بعد از تعداد مشخصی از حرکات اپدیت شده و ورن دهی میکند تا در نهایت بعد از تعدادی تکرار مقادیر مربوط به هر درابه به مقداری همگرا می شوند. در الگوریتم Deep-Q-Learning در حالتی که فضای سوال بزرگ باشد و اپدیت کردن ماتریس Q کاری بسار دشوار خواهد بود. در این سوال نیز فضای سیستم بسیار گسترده است و تعداد حالات زیادی باید مورد بررسی قرار گیرد به همین دلیل از این روش استفاده می شود. در این قسمت به بررسی قسمت هایی از کد که مربوط به این متد می باشد میپردازیم:



شکل ۲: مقایسه دو روش

کلاس ExperienceRelay تعریف شده برای ذخیره سازی ۴ متغیر می باشد که این مقادیر در هر اپیزود تکرار میشود . این متغیر ها برابر هستند با :

۱. حالت کنونی

۲. لحظه بعد

۳. اکشن های موجود در این حالت

۴. پاداش هر حرکت

این مقادیر به صورت یک رشته tuple در کنار یکدیگر قرار میگیرند و برای هر خانه و هر state این مقادیر به دست می آیند.

$$Experience_i = (state_i, action_i, state_{i+1}, reward)$$

و همان طوز که ذکر شد این مقادیر ردر کنار یکدیگر قرار میگیرند تا برای تمامی خانه های سوال این tuple ها شکل گیرند.

$$Total_{experince} = [e_1, e_2, e_3, ..., e_n]$$

حال در قسمت آموزش mini-batch هایی را به صورت رندوم انتخاب میکنیم. این نمونه ها به صورت تجربه هایی در محیط هستند. در این کلاس سه متد وجود دارد.

۱. init : در این روش از یک DEQ(Double Ended Queue) استفاده میکنیم برای ذخیره تجربیات

۲. store-trans : تنها برای ذخیره سازی اطلاعات در حافظه

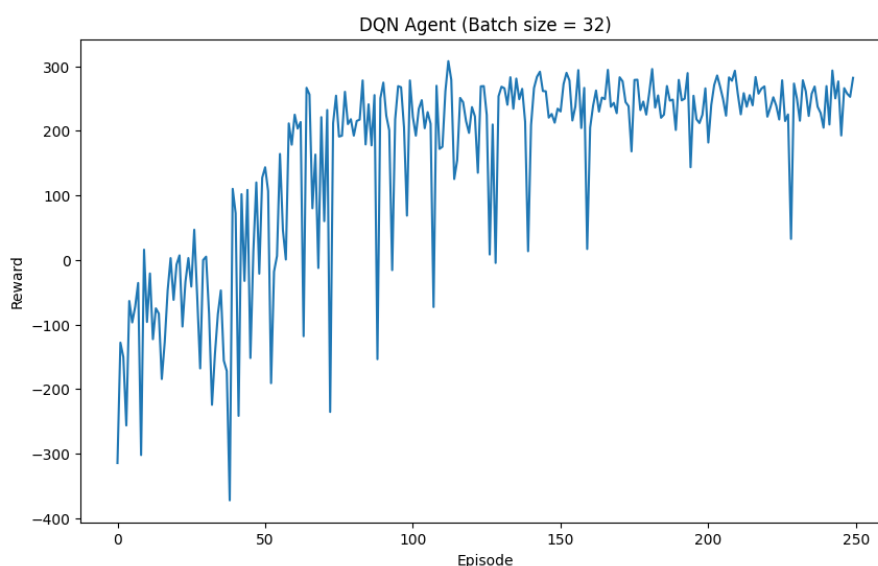
۳. Sample : نمونه برداری از حافظه به صورت تصادفی به تعداد Batch Size

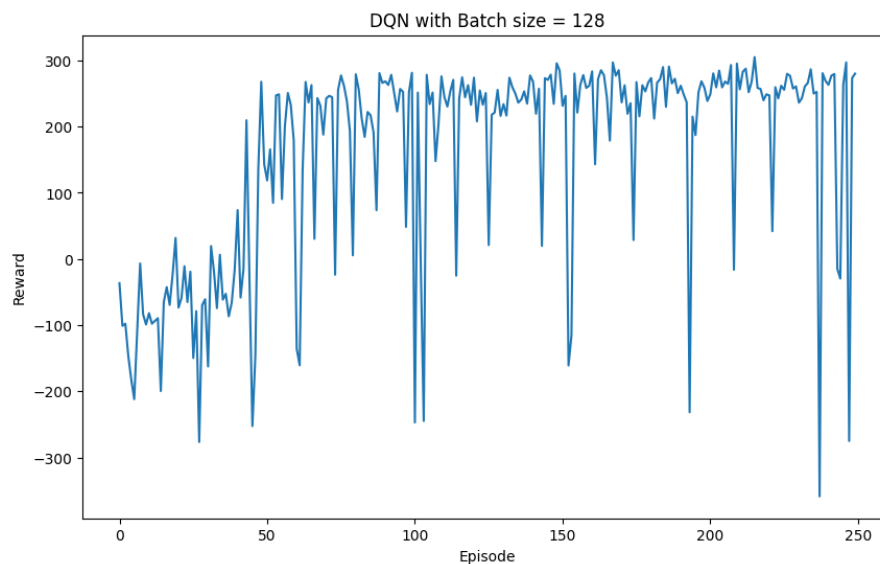
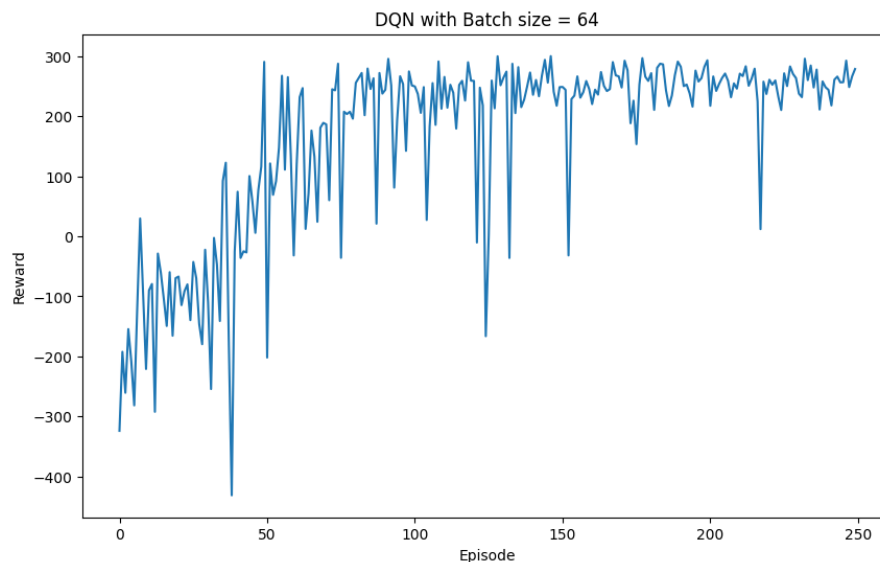
کلاس DeepQNetwork مربوط به ساختار شبکه عصبی است. در این شبکه تعداد لایه ورودی برابر با ۸ می باشد که برابر با اندازه فضای حالت می باشد. و از طرفی هم تعداد لایه های خروجی برابر با ۴ که هماتد تعداد اکشن است می باشد. تعداد نورون های هر لایه برابر با ۵۱۲ تا است. در تمامی لایه های میانی (پنهان) خروجی تمامی لایه ها نرمالایز می شوند با استفاده از Layernorm و نیز تمامی تابع فعال ساز این لایه ها ReLU می باشد. برای جلوگیری از overfitting از روش drop out استفاده می شود.

کلاس DQNAgent : در این کلاس ابتدا هاپیر پارامتر ها، ساختار شبکه و نوع بهینه ساز با استفاده از متد init تعیین می شوند. سپس برای اینکه بازیکن بتواند اقدام بعدی را انتخاب کند از استراتژی epsilon greedy استفاده میکنیم. تمام مراحل مربوط به انتخاب اکشن بعدی در take action انتخاب می شوند. در قسمت update-params همان طور که از اسم این بخش مشخص است برای اپدیت کردن پارامتر ها در هر مرحله استفاده میشود سپس این اطلاعات به روز رسانی شده را با استفاده از load , save ذخیره میکنیم.

حال در بخش آموزش میبایست مراحل زیر انجام شود. ابتدا می بایست مقادیر پارامتر های آموزش تعیین شوند و سپس فرایند اصلی آغاز میشود. یعنی با توجه به موقعیتی که سفینه در آن قرار دارد یک اکشن انتخاب میشود و در قسمت Experience ذخیره میشود سپس با توجه با اقدام انتخاب شده وانجام شده تجربیات مربوط به آن ذخیره میشود. در ضمن در ابتدا مقدار epsilon زیاد می باشد در واقع در ابتدا سیستم در حالت جستجوی در کل قسمت جواب می باشد و سپس با انجام دادن اقداماتی به تدریج این مقدار کاهش میابد تا به مسیر بهینه دست بیابیم. در واقع یک تعادل بین جستجو در کل مسیر ها و مسیر بهینه میباشد.

در ادامه پاداش تجمعی برای Batch size های مختلف ترسیم شده است:





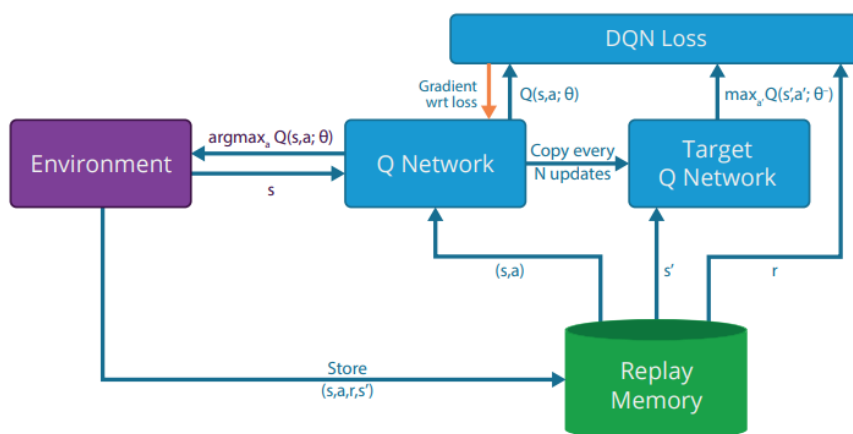
در هر سه نمودار روند مشابهی قابل رویت است . در واقع با افزایش اپیزود ها مقدار پاداش افزایش میابد. در واقع این به معنای بهبود عملکرد سیستم می باش پرا که از انجام اکشن هایی امتیاز منفی دارند امتناع کرده است. حال باید بهترین batch size را انتخاب کنیم از روی روند کلی و شمای ان نمیتون به این امر دست یافت . پس در چندین اپیزود پابان به مقایسه می پردازیم . دیده میشود که در حالت ۱۲۸ در حالات پایانی مقدار نوسانات به نسبت بقیه بیشتر می باشد و نیز مقدار ان تیز برابر با -۳۰۰ میباشد که گویا بعد از تکرار های مکرر عملکرد خوبی از خود نشان داده است. حال به سراغ ۶۴ میرویم . برای این مقدار نیز دز اپیزود های پایانی نوسان داریم ولی مقدار انها کمتر و شدتشان بیشتر میباشد. در حالت ۳۲ تمامی این موارد بهبود یافته است طوری که شدت و مقدار نوسانات به شدت کاهش یافته است. پس در نهایت مقدار ۳۲ انتخاب میشود. دقت شود معیار پشیمانی برابر با اختلاف میان حرکت در حال انجام با بهترین حالت ممکنه می باشد. در واقع در این سوال بهترین مقدار برابر با امتیاز ۳۰۰ میباشد و هر چه قدر این اختلاف بیشتر باشد و منفی تر میزان پشیمانی بیشتر است. در مورد نمودار های پیوسته این مقدار برابر با مساحت بالای نمودار تا خط امتیاز ۳۰۰ می باشد. که با توجه به این توضیحات ابن مساحت برای حالت ۶۴ در کمترین حالت خود قرار دارد و معیار پشیمانی برای ۶۴ کمتر است. برای مشاهده فیلم عملکرد عامل در اپیزودهای مختلف به ازای دسته ۳۲ نمونه ای .

شماره اپیزود	لینک
۵۰	لینک
۱۰۰	لینک
۱۵۰	لینک
۲۰۰	لینک
۲۵۰	لینک

جدول ۱: فیلم های عملکرد عامل در روش DQN

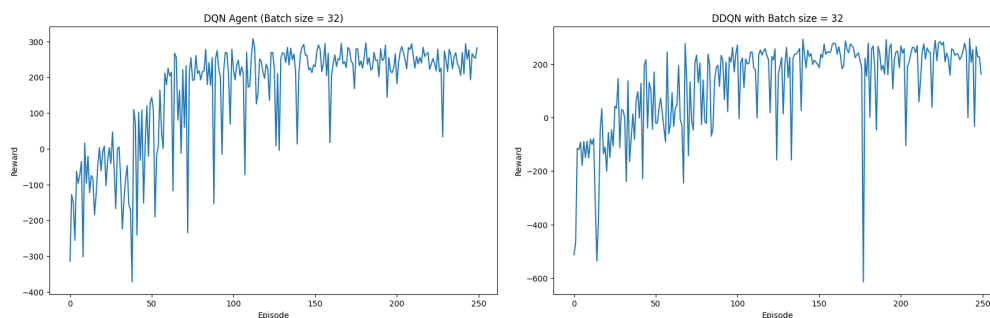
### ۳.۱ بخش سوم

در ابتدای سوال در مورد DQN , DDQN صحبت شد حال به بررسی بستر تفاوت این دو میپردازیم. در روش DDQN از دو شبکه عصبی استفاده است به نام های target network و evaluate network . شبکه دوم اضافه شده جهت افزایش دقت در محاسبه ماتریس Q این امر سبب افزایش سرعت همگرایی سیستم و پایداری بیشتر آن میشود. در ضمن در این روش حجم محاسبات و پردازش بیشتر می باشد و از طرفی نیز احتمال رسیدن به بهینه ترین حالت بیشتر است زیرا که تمامی موارد مورد بررسی قرار میگیرند.



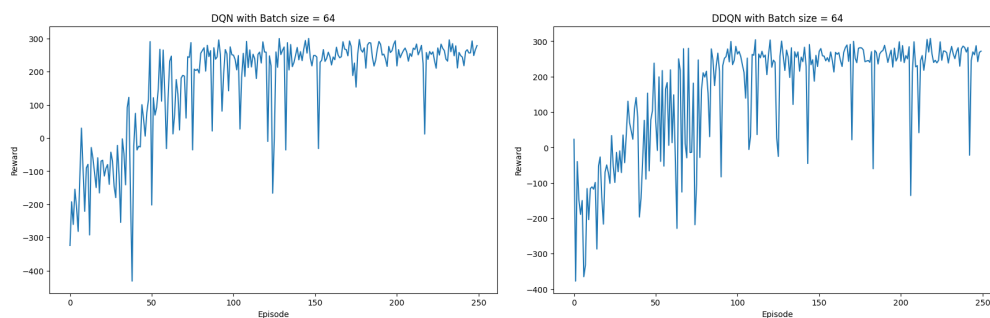
شکل ۳: ساختار DDQN

حال به بررسی کد داده شده میپردازیم. بسیاری از قسمت های کد همانند حالت DQN می باشد. دقت شود در کلاس DDQNAgent تفاوت هایی دارد نظیر اینکه دو شبکه با نام های value net , target net قرار دارند که ساختار کلی آنها همانند یکدیگر است. در ضمن تابع هزینه استفاده شده در این قسمت MSE حداقل مربعات است. حال نمودار های مربوط به عملکرد های مختلف را رسم میکنیم.

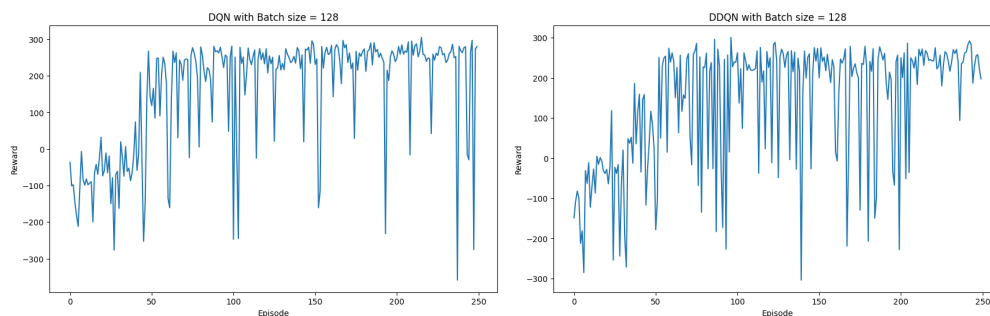


شکل ۴: دسته های ۳۲ عددی





شکل ۵: دسته های ۶۴ عددی



شکل ۶: دسته های ۱۲۸ عددی

با مقایسه نمودار ها دیده میشود که نمودار مربوط به تمامی مقداری عملکردشان بهبود یافته و شدت و میزان نوسانات در پایان اپیزود ها کاهش یافته است و البته دقت شود در این حالت بهترین مقدار برای ۱۲۸ می باشد چرا که کمترین مقدار نوسانات و کمترین شدت نوسانات را دارد.. فیلم مربوط به عملکرد سیستم در اپیزودهای ۱۰۰ و ۲۵۰ به ازای الگوریتم های مختلف و batch size های مختلف

الگوریتم	شماره اپیزود	Mini Batch	لینک
DDQN	۱۰۰	۳۲	<a href="#">لینک</a>
DQN	۱۰۰	۳۲	<a href="#">لینک</a>
DDQN	۲۵۰	۳۲	<a href="#">لینک</a>
DQN	۲۵۰	۳۲	<a href="#">لینک</a>
DDQN	۱۰۰	۶۴	<a href="#">لینک</a>
DQN	۱۰۰	۶۴	<a href="#">لینک</a>
DDQN	۲۵۰	۶۴	<a href="#">لینک</a>
DQN	۲۵۰	۶۴	<a href="#">لینک</a>
DDQN	۱۰۰	۱۲۸	<a href="#">لینک</a>
DQN	۱۰۰	۱۲۸	<a href="#">لینک</a>
DDQN	۲۵۰	۱۲۸	<a href="#">لینک</a>
DQN	۲۵۰	۱۲۸	<a href="#">لینک</a>

```

1 # Install packages
2 !apt-get update
3 !apt-get install -y swig
4 !sudo apt-get update
5 !sudo apt-get install xvfb
6 !pip install rarfile --quiet
7 !pip install stable-baselines3 > /dev/null
8 !pip install box2d-py > /dev/null
9 !pip install gym pyvirtualdisplay > /dev/null 2>&1
10 !sudo apt-get install -y xvfb python-opengl ffmpeg > /dev/null 2>&1
11 import gym
12 import io
13 import os
14 import glob
15 import torch
16 import base64
17 import numpy as np
18 import copy
19 import matplotlib.pyplot as plt
20 from stable_baselines3 import DQN
21 from stable_baselines3.common.results_plotter import ts2xy, load_results
22 from stable_baselines3.common.callbacks import EvalCallback
23 from gym.wrappers import RecordVideo
24 from IPython.display import HTML
25 from IPython import display as ipythondisplay
26 from pyvirtualdisplay import Display
27 import random
28 from collections import namedtuple, deque
29 import torch.nn as nn
30 import torch.nn.functional as F
31 import torch.optim as optim
32
33 # Creat envoiroment
34 env = gym.make("LunarLander-v2", new_step_api=True)
35
36 env.reset()
37
38 state_size = env.observation_space.shape[0]
39 action_size = env.action_space.n
40
41 # Introduction to state and action space
42 print(f'size state space : {state_size}')
43 print(f'size of action space: {action_size}')
44
45 print(f'Sample state space: {env.observation_space.sample()}')
46 print(f'Sample action space: {env.action_space.sample()}')
47
48 # Set up virtual display
49 display = Display(visible=0, size=(1400, 900))
50 display.start()
51
52 # Utility function to enable video recording of gym environment and displaying
  it
53 def show_video():
54     mp4list = glob.glob('video/*.mp4')
55     if len(mp4list) > 0:

```

```

56     mp4 = mp4list[0]
57     video = io.open(mp4, 'r+b').read()
58     encoded = base64.b64encode(video)
59     ipythondisplay.display(HTML(data='<video alt="test" autoplay loop
controls style="height: 400px;">
60         <source src="data:video/mp4;base64,{0}" type="video/mp4" />
61         </video>'''.format(encoded.decode('ascii'))))
62     else:
63         print("Could not find video")
64
65     # Set random seed
66 SEED = 14
67 random.seed(SEED)
68 torch.manual_seed(SEED)
69
70 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
71 device
72
73 # Experience replay
74 Transition = namedtuple('Transition',
75                         ('state', 'action', 'next_state', 'reward', 'done'))
76
77 class ExperienceReplay():
78     def __init__(self, capacity) -> None:
79         self.memory = deque([], maxlen=capacity)
80
81     def store_trans(self, s, a, sp, r, done):
82         transition = Transition(s, a, sp, r, done)
83         self.memory.append(transition)
84
85     def sample(self, batch_size):
86         return random.sample(self.memory, batch_size)
87
88     def __len__(self):
89         return len(self.memory)
90
91 # Deep Q-Network Structure
92
93 class DeepQNetwork(nn.Module):
94     def __init__(self, state_size, action_size) -> None:
95         super(DeepQNetwork, self).__init__()
96         net_list = nn.ModuleList([
97             torch.nn.Linear(state_size, 512),
98             torch.nn.ReLU(),
99             torch.nn.LayerNorm(512),
100             torch.nn.Dropout(0.1),
101             torch.nn.Linear(512, 512),
102             torch.nn.ReLU(),
103             torch.nn.LayerNorm(512),
104             torch.nn.Dropout(0.1),
105             torch.nn.Linear(512, 512),
106             torch.nn.ReLU(),
107             torch.nn.Linear(512, action_size)
108         ])
109         self.net = torch.nn.Sequential(*net_list).to(device)
110
111     def forward(self, x):
112         x = x.to(device)
113         x = self.net(x)
114         return x

```

```

115
116 # DQN agent
117
118 class DQNAgent():
119     def __init__(self, state_size, action_size, batch_size,
120                 gamma=0.99, buffer_size=25000, alpha=1e-4):
121         self.state_size = state_size
122         self.action_size = action_size
123         self.batch_size = batch_size
124         self.gamma = gamma
125         self.experience_replay = ExperienceReplay(buffer_size)
126         self.value_net = DeepQNetwork(state_size, action_size).to(device)
127         self.optimizer = optim.Adam(self.value_net.parameters(), lr=alpha)
128
129     def take_action(self, state, eps=0.0):
130         self.value_net.eval()
131         if random.random() > eps:
132             with torch.no_grad():
133                 state_tensor = torch.tensor(state, dtype=torch.float32).
134                 unsqueeze(0).to(device)
135                 action_values = self.value_net(state_tensor)
136                 return torch.argmax(action_values).item()
137         else:
138             return np.random.randint(0, self.action_size)
139
140     def update_params(self):
141         if len(self.experience_replay) < self.batch_size:
142             return
143         batch = Transition(*zip(*self.experience_replay.sample(self.batch_size)
144 ))
145         state_batch = torch.tensor(batch.state, dtype=torch.float32).to(device)
146         action_batch = torch.tensor(batch.action).unsqueeze(1).to(device)
147         next_state_batch = torch.tensor(batch.next_state, dtype=torch.float32).
148         to(device)
149         reward_batch = torch.tensor(batch.reward, dtype=torch.float32).unsqueeze
150         (1).to(device)
151         done_batch = torch.tensor(batch.done, dtype=torch.float32).unsqueeze(1).
152         to(device)
153
154         self.value_net.train()
155         q_expected = self.value_net(state_batch).gather(1, action_batch)
156         q_targets_next = self.value_net(next_state_batch).detach().max(1)[0].
157         unsqueeze(1)
158         q_targets = reward_batch + (self.gamma * q_targets_next * (1 -
159         done_batch))
160         loss = nn.functional.mse_loss(q_expected, q_targets)
161
162         self.optimizer.zero_grad()
163         loss.backward()
164         self.optimizer.step()
165
166     def save(self, fname):
167         torch.save(self.value_net.state_dict(), fname)
168
169     def load(self, fname, device):
170         self.value_net.load_state_dict(torch.load(fname, map_location=device))
171
172 n_episodes = 250
173 eps = 1.0
174 eps_decay_rate = 0.97
175 eps_end = 0.01

```

```

168 BATCH_SIZE = 128
169
170
171
172 # TODO: create agent
173 agent = DQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
174
175 crs = np.zeros(n_episodes) # cummulative rewards
176 crs_recent = deque(maxlen=25) # recent cummulative rewards
177
178 # training loop
179 for i_episode in range(1, n_episodes+1):
180     # TODO: initialize the environment and state
181     if i_episode % 50 == 0:
182         env = RecordVideo(gym.make("LunarLander-v2"), f"./DQN/batch{BATCH_SIZE}/
eps{i_episode}")
183     else:
184         env = gym.make("LunarLander-v2")
185         state = env.reset()
186         done = False
187         cr = 0 # episode cummulative rewards
188         while not done:
189             env.render()
190             # TODO: select and perform an action
191             action = agent.take_action(state, eps)
192             #
193             print(action)
194             next_state, reward, done, info = env.step(action)
195             # TODO: store transition in experience replay
196             agent.experience_replay.store_trans(state, action, next_state, reward,
done)
197             # TODO: update agent
198             agent.update_params()
199             # TODO: update current state and episode cummulative rewards
200             state = next_state
201             cr += reward
202
203             # TODO: decay epsilon
204             eps = eps * eps_decay_rate
205             eps = max(eps, eps_end)
206             # TODO: update current cummulative rewards and recent cummulative rewards
207             crs[i_episode - 1] = cr
208             crs_recent.append(cr)
209             # TODO: save agent every 50 episodes
210             if i_episode % 50 == 0:
211                 agent.save(f"DQN_batch{BATCH_SIZE}_eps{i_episode}.pt")
212
213             # print logs
214             print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(
i_episode, np.mean(crs_recent), eps), end="")
215             if i_episode % 25 == 0:
216                 print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(
i_episode, np.mean(crs_recent), eps))
217 # Plot cumulative reward of each episode
218 fig = plt.figure()
219 plt.figure(figsize=(10,6))
220 plt.plot(np.arange(len(crs)), crs)
221 plt.ylabel('Reward')
222 plt.xlabel('Episode')
223 plt.title(f"DQN with Batch size = {BATCH_SIZE}")
224 plt.show()

```

```

224
225 # Double DQN Agent
226
227 class DDQNAgent():
228     # NOTE: DON'T change initial values
229     def __init__(self, state_size, action_size, batch_size,
230                 gamma=0.99, buffer_size=25000, alpha=1e-4):
231         # network parameter
232         self.state_size = state_size
233         self.action_size = action_size
234
235         # hyperparameters
236         self.batch_size = batch_size
237         self.gamma = gamma
238
239         # experience replay
240         self.experience_buffer = ExperienceReplay(buffer_size)
241
242         # networks
243         self.value_net = DeepQNetwork(state_size, action_size).to(device)
244         self.target_value_net = DeepQNetwork(state_size, action_size).to(device)
245         self.update_target_network()
246
247         # optimizer
248         # TODO: create adam for optimizing network's parameter (learning rate=
alpha)
249         # NOTE: target network parameters DOSEN'T update with optimizer
250         self.optimizer = optim.Adam(self.value_net.parameters(), lr=alpha)
251
252     def take_action(self, state, eps=0.0):
253         # TODO: take action using e-greedy policy
254         # NOTE: takes action using the greedy policy with a probability of 1
and a random action with a probability of
255         # NOTE:
256         self.value_net.eval()
257         rand_eps = random.random()
258         if rand_eps > eps:
259             with torch.no_grad():
260                 return torch.argmax(self.value_net(torch.tensor(state).to(device)
))) .detach().cpu().numpy()
261         else:
262             return np.random.randint(0, self.action_size)
263
264     def update_params(self):
265         if len(self.experience_buffer) < self.batch_size:
266             return
267         # transition batch
268         batch = Transition(*zip(*self.experience_buffer.sample(self.batch_size)
))
269
270         state_batch = torch.from_numpy(np.vstack(batch.state)).float().to(device)
271
272         action_batch = torch.tensor(np.vstack(batch.action)).long().to(device)
273         next_state_batch = torch.from_numpy(np.vstack(batch.next_state)).float().to(device)
274         reward_batch = torch.tensor(np.vstack(batch.reward)).float().to(device)
275         done_batch = torch.tensor(np.vstack(batch.done)).to(device)
276
277         # calculate loss w.r.t DQN algorithm
278         self.value_net.train()

```

```

278         # STEP1
279         q_targets_next = self.target_value_net(next_state_batch).detach().max(1)
280         [0].unsqueeze(1)
281         # STEP2
282         # TODO: compute Q values [Q(s_t, a)]
283         q_targets = reward_batch + self.gamma * q_targets_next * (1 - done_batch
284         *1)
285         q_expected = self.value_net(state_batch).gather(1, action_batch)
286         # STEP3
287         # TODO: compute mse loss
288         loss = nn.functional.mse_loss(q_expected, q_targets)
289         # TODO: optimize the model
290         # NOTE: DON'T forget to set the gradients to zeros
291         self.optimizer.zero_grad()
292         loss.backward()
293         self.optimizer.step()
294
295     def update_target_network(self):
296         # TODO: copy main network parameters to target network parameters
297         self.target_value_net = copy.deepcopy(self.value_net)
298
299     def save(self, fname):
300         # TODO: save checkpoint
301         torch.save(self.value_net, fname)
302
303     def load(self, fname, device):
304         # TODO: load checkpoint
305         self.value_net = torch.load(fname, device)
306
307 # Set training parameters
308 n_episodes = 250
309 eps = 1.0
310 eps_decay_rate = 0.97
311 eps_end = 0.01
312 BATCH_SIZE = 32
313
314 # Training phase
315
316 # TODO: create agent
317 agent = DDQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
318
319 crs = np.zeros(n_episodes) # cummulative rewards
320 crs_recent = deque(maxlen=25) # recent cummulative rewards
321
322 # training loop
323 for i_episode in range(1, n_episodes+1):
324     # TODO: initialize the environment and state
325     if i_episode % 50 == 0:
326         env = RecordVideo(gym.make("LunarLander-v2"), f"./DDQN/batch{BATCH_SIZE}
327         {i_episode}")
328     else:
329         env = gym.make("LunarLander-v2")
330         state = env.reset()
331         done = False
332         cr = 0 # episode cummulative rewards
333         action_count = 0
334         while not done:
335             env.render()
336             # TODO: select and perform an action
337             action = agent.take_action(state, eps)

```

```

335     next_state, reward, done, info = env.step(action)
336     # TODO: store transition in experience replay
337     agent.experience_buffer.store_trans(state, action, next_state, reward,
done)
338     # TODO: update agent
339     agent.update_params()
340     # TODO: update current state and episode cumulative rewards
341     state = next_state
342     cr += reward
343     action_count += 1
344     if action_count % 5 == 0:
345         agent.update_target_network()
346
347     # TODO: decay epsilon
348     eps = eps * eps_decay_rate
349     eps = max(eps, eps_end)
350
351     # TODO: update current cumulative rewards and recent cumulative rewards
352     crs[i_episode - 1] = cr
353     crs_recent.append(cr)
354
355     # TODO: save agent every 50 episodes
356     if i_episode % 50 == 0:
357         agent.save(f"DDQN/DDQN_batch{BATCH_SIZE}_eps{i_episode}.pt")
358
359     # print logs
360     print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(
i_episode, np.mean(crs_recent), eps), end="")
361     if i_episode % 25 == 0:
362         print('\rEpisode {} \tAverage Reward: {:.2f} \tEpsilon: {:.2f}'.format(
i_episode, np.mean(crs_recent), eps))
363
364 fig = plt.figure()
365 plt.figure(figsize=(10,6))
366 plt.plot(np.arange(len(crs)), crs)
367 plt.ylabel('Reward')
368 plt.xlabel('Episode')
369 plt.title(f"DDQN with Batch size = {BATCH_SIZE}")
370 plt.show()

```

Listing 1: Code Question 1