



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

درس یادگیری ماشین گزارش مینی پروژه ۴

نام و نام خانوادگی	دانیال عبداللہی نژاد
شماره دانشجویی	۴۰۱۰۹۱۶۴
تاریخ	تیرماه ۱۴۰۳



فهرست مطالب

۳	۱ سوال اول
۳	۲ سوال دوم
۱۳	۱.۲
۱۵	۲.۲
۱۸	۳.۲



فهرست تصاویر

۱۳ Lunar lander محیط	۱
۱۶ batch size=32 نمودار reward برای	۲
۱۶ batch size=64 نمودار reward برای	۳
۱۷ batch size=128 نمودار reward برای	۴
۲۰ batch size=64 نمودار reward برای	۵
۲۰ batch size=128 نمودار reward برای	۶
۲۱ نمونه خروجی ویدیویی	۷



۱ سوال اول

۲ سوال دوم

در ابتدای کتابخانه‌ها به صورت زیر برای استفاده در مراحل بعد وارد می‌شوند.

```
1 import gym
2 import io
3 import os
4 import glob
5 import torch
6 import base64
7 import numpy as np
8 import matplotlib.pyplot as plt
9 from stable_baselines3 import DQN
10 from stable_baselines3.common.results_plotter import ts2xy, load_results
11 from stable_baselines3.common.callbacks import EvalCallback
12 from gym.wrappers import RecordVideo
13 from IPython.display import HTML
14 from IPython import display as ipythondisplay
15 from pyvirtualdisplay import Display
16 import random
17 from collections import namedtuple, deque
18 import torch.nn as nn
19 import torch.nn.functional as F
20 import torch.optim as optim
21
22 display = Display(visible=0, size=(1400, 900))
23 display.start()
```

کتابخانه stable-baselines3 یک مجموعه پیاده‌سازی‌های پایدار و بهینه از الگوریتم‌های یادگیری تقویتی در پایتون است. این کتابخانه به راحتی قابل استفاده و توسعه است و شامل الگوریتم‌های معروفی مانند DQN، PPO، A2C و SAC می‌باشد.

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



این تابع برای نمایش ویدیوهای ضبط شده از محیط Gym استفاده می‌شود. ویدیوها به صورت base64 رمزگذاری شده و در مرورگر نمایش داده می‌شوند.

```
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)
```

Experience Replay یکی از اجزای کلیدی در پیاده‌سازی‌های DQN و DDQN در کد ارائه شده است. این بخش به عامل اجازه می‌دهد تا تجربیات خود را در یک حافظه ذخیره کند و از آن‌ها به صورت نمونه‌گیری تصادفی در فرآیند آموزش استفاده کند. Experience Replay با استفاده از یک صف دوتایی (deque) پیاده‌سازی شده است که ظرفیت مشخصی دارد. این حافظه تجربیات عامل را به صورت transition شامل state، action، next state، reward و done ذخیره می‌کند. متد store_trans برای ذخیره‌سازی transition در حافظه استفاده می‌شود. این متد یک transition جدید شامل موارد گفته شده را به حافظه اضافه می‌کند.

```
1 def store_trans(self, s, a, sp, r, done):
2     transition = Transition(s, a, sp, r, done)
3     self.memory.append(transition)
```

sample برای نمونه‌گیری تصادفی از transition در حافظه استفاده می‌شود. این متد تعداد مشخصی از transition‌ها را به صورت تصادفی از حافظه انتخاب می‌کند تا در فرآیند به‌روزرسانی شبکه عصبی استفاده شوند.

```
1 def sample(self, batch_size):
2     return random.sample(self.memory, batch_size)
```

در فرآیند آموزش، عامل transition را در حافظه Experience Replay ذخیره می‌کند و سپس به صورت دوره‌ای از این حافظه نمونه‌گیری کرده و شبکه عصبی خود را به‌روزرسانی می‌کند. این روش به کاهش همبستگی بین داده‌های ورودی کمک می‌کند و به پایداری و کارایی بهتر مدل منجر می‌شود.

```
1 for i_episode in range(1, n_episodes+1):
2     ...
3     while not done:
4         ...
5         agent.experience_replay.store_trans(state, action, next_state, reward, done or truncated)
6         agent.update_params()
7         ...
```



این تکنیک از بیش‌برازش به داده‌های اخیر جلوگیری می‌کند.

```

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)

```

کلاس DeepQNetwork با استفاده از یک شبکه عصبی عمیق، تخمینی از Q-values را برای هر اقدام در یک حالت خاص تولید می‌کند که به عامل کمک می‌کند تا بهترین اقدام ممکن را انتخاب کند.

ورودی این شبکه، حالت (state) محیط است که به صورت یک بردار با اندازه state_size ارائه می‌شود. این بردار می‌تواند شامل ویژگی‌های مختلفی از محیط باشد که برای تصمیم‌گیری عامل مهم هستند. به عنوان مثال، در محیط LunarLander-v2، ورودی شامل موقعیت و سرعت فضاپیما و زاویه و نرخ چرخش آن است.

خروجی شبکه، یک بردار با اندازه action_size است که هر عنصر آن نشان‌دهنده Q-value برای یک اقدام خاص است. این Q-values تخمینی از مقدار کل پاداشی است که عامل می‌تواند از هر حالت و با انجام هر اقدام خاص انتظار داشته باشد. شبکه DeepQNetwork شامل لایه‌های زیر است:

- **لایه Linear:** این لایه‌ها عملیات خطی را بر روی ورودی انجام می‌دهند و هر کدام شامل تعداد مشخصی نورون هستند. در این شبکه، سه لایه خطی با ۵۱۲ نورون وجود دارد.
- **لایه ReLU:** این لایه‌ها یک تابع فعال‌سازی غیرخطی را اعمال می‌کنند که به مدل اجازه می‌دهد تا روابط غیرخطی پیچیده را یاد بگیرد. ReLU مخفف Rectified Linear Unit است.
- **لایه LayerNorm:** این لایه‌ها نرمال‌سازی لایه‌ای را انجام می‌دهند که به پایداری و همگرایی سریع‌تر آموزش کمک می‌کند.
- **لایه Dropout:** این لایه‌ها به منظور جلوگیری از بیش‌برازش استفاده می‌شوند و برخی از نورون‌ها را به صورت تصادفی در هر دوره آموزشی غیر فعال می‌کنند.

تابع forward برای عبور ورودی از شبکه و تولید خروجی استفاده می‌شود. این تابع ورودی x را از طریق تمام لایه‌های تعریف شده در __init__ عبور می‌دهد و خروجی نهایی را که شامل Q-values است، باز می‌گرداند.



```
1 def forward(self, x):  
2     return self.net(x)
```

کلاس DQNAgent نماینده عامل DQN است. این کلاس شامل توابعی برای انتخاب اقدام، به‌روزرسانی پارامترها و ذخیره و بارگذاری مدل است. عامل از یک شبکه عصبی برای انتخاب و ارزیابی اقدامات استفاده می‌کند.

کلاس DQNAgent نماینده عامل Lunar Lander است. در هر اپیزود، این عامل با محیط تعامل می‌کند تا بهترین اقدامات را بر اساس حالت‌های فعلی انجام دهد.

```
1 class DQNAgent():  
2     def __init__(self, state_size, action_size, batch_size, gamma=0.99, buffer_size=25000, alpha  
3         =1e-4):  
4         self.state_size = state_size  
5         self.action_size = action_size  
6         self.batch_size = batch_size  
7         self.gamma = gamma  
8         self.experience_replay = ExperienceReplay(buffer_size)  
9         self.value_net = DeepQNetwork(state_size, action_size).to(device)  
10        self.optimizer = optim.Adam(self.value_net.parameters(), lr=alpha)  
11  
12    def take_action(self, state, eps=0.0):  
13        self.value_net.eval()  
14        if random.random() > eps:  
15            with torch.no_grad():  
16                return torch.argmax(self.value_net(torch.tensor(state).float().unsqueeze(0).to(  
17                    device))).item()  
18            else:  
19                return np.random.randint(0, self.action_size)  
20  
21    def update_params(self):  
22        if len(self.experience_replay) < self.batch_size:  
23            return  
24        batch = Transition(*zip(*self.experience_replay.sample(self.batch_size)))  
25  
26        state_batch = torch.tensor(np.array(batch.state), dtype=torch.float32).to(device)  
27        action_batch = torch.tensor(batch.action, dtype=torch.int64).unsqueeze(1).to(device)  
28        next_state_batch = torch.tensor(np.array(batch.next_state), dtype=torch.float32).to(  
29            device)  
30        reward_batch = torch.tensor(batch.reward, dtype=torch.float32).unsqueeze(1).to(device)  
31        done_batch = torch.tensor(batch.done, dtype=torch.float32).unsqueeze(1).to(device)  
32  
33        q_expected = self.value_net(state_batch).gather(1, action_batch)  
34        q_targets_next = self.value_net(next_state_batch).detach().max(1)[0].unsqueeze(1)  
35        q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch))  
36        loss = F.mse_loss(q_expected, q_targets)  
37  
38        self.optimizer.zero_grad()
```



```

36     loss.backward()
37     self.optimizer.step()
38
39     def save(self, fname):
40         torch.save(self.value_net.state_dict(), fname)
41
42     def load(self, fname):
43         self.value_net.load_state_dict(torch.load(fname, map_location=device))

```

در هر اپیزود، عامل با تنظیمات اولیه محیط شروع می‌کند. این شامل reset محیط و تنظیم حالت اولیه است. این تابع مقداردهی اولیه برای کلاس DQNAgent را انجام می‌دهد و پارامترهای مختلف را تنظیم می‌کند:

- **state_size**: اندازه بردار حالت‌ها.
- **action_size**: تعداد اقدامات ممکن.
- **batch_size**: اندازه دسته‌ها برای نمونه‌گیری از حافظه Experience Replay.
- **gamma**: نرخ تخفیف برای محاسبه پاداش‌های آینده.
- **buffer_size**: ظرفیت حافظه Experience Replay.
- **alpha**: نرخ یادگیری برای به‌روزرسانی پارامترهای شبکه.

```

1 def __init__(self, state_size, action_size, batch_size, gamma=0.99, buffer_size=25000, alpha=1e
  -4):
2     self.state_size = state_size
3     self.action_size = action_size
4     self.batch_size = batch_size
5     self.gamma = gamma
6     self.experience_replay = ExperienceReplay(buffer_size)
7     self.value_net = DeepQNetwork(state_size, action_size).to(device)
8     self.optimizer = optim.Adam(self.value_net.parameters(), lr=alpha)

```

در هر گام زمانی، عامل یک اقدام را با استفاده از سیاست e-greedy انتخاب می‌کند. اگر عدد تصادفی کمتر از eps باشد، اقدام تصادفی انتخاب می‌شود (exploration). در غیر این صورت، اقدام با بالاترین Q-value تخمین زده شده توسط شبکه عصبی انتخاب می‌شود (exploitation).

عامل اقدام انتخاب شده را اجرا می‌کند و محیط حالت جدید، پاداش، و وضعیت اتمام (done) را باز می‌گرداند.

```

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(device)
6             )).item()
7     else:
8         return np.random.randint(0, self.action_size)

```




این تابع پارامترهای شبکه عصبی را به‌روزرسانی می‌کند:

- اگر تعداد تجربیات در حافظه کمتر از `batch_size` باشد، از به‌روزرسانی خودداری می‌کند.
- انتقال‌ها را به صورت دسته‌ای نمونه‌گیری می‌کند و بردارهای حالت، اقدام، حالت بعدی، پاداش و علامت پایان را ایجاد می‌کند.
- `Q-value` های فعلی و هدف را محاسبه می‌کند و از `loss` میانگین مربعات خطا (MSE) برای به‌روزرسانی شبکه استفاده می‌کند.

```
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(device)
9     reward_batch = torch.tensor(batch.reward, dtype=torch.float32).unsqueeze(1).to(device)
10    done_batch = torch.tensor(batch.done, dtype=torch.float32).unsqueeze(1).to(device)
11
12    q_expected = self.value_net(state_batch).gather(1, action_batch)
13    q_targets_next = self.value_net(next_state_batch).detach().max(1)[0].unsqueeze(1)
14    q_targets = reward_batch + (self.gamma * q_targets_next * (1 - done_batch))
15    loss = F.mse_loss(q_expected, q_targets)
16
17    self.optimizer.zero_grad()
18    loss.backward()
19    self.optimizer.step()
```

این تابع مدل یادگیری شده را ذخیره می‌کند.

```
1 def save(self, fname):
2     torch.save(self.value_net.state_dict(), fname)
```

این تابع مدل یادگیری شده را بارگذاری می‌کند.

```
1 def load(self, fname):
2     self.value_net.load_state_dict(torch.load(fname, map_location=device))
```

این فرآیند در هر اپیزود تکرار می‌شود تا عامل بتواند بهترین استراتژی را برای تعامل با محیط Lunar Lander یاد بگیرد.

حلقه آموزش برای اندازه دسته ۳۲

```
1 env = gym.make('LunarLander-v2', render_mode="rgb_array")
2 state_size = env.observation_space.shape[0]
3 action_size = env.action_space.n
4 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
5
```



```
6 n_episodes = 250
7 eps = 1.0
8 eps_decay_rate = 0.97
9 eps_end = 0.01
10 BATCH_SIZE = 32
11
12 agent = DQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
13 crs = np.zeros(n_episodes)
14 crs_recent = deque(maxlen=25)
15
16 for i_episode in range(1, n_episodes+1):
17     env = RecordVideo(gym.make("LunarLander-v2", render_mode="rgb_array"), f"./DQN/batch{
18         BATCH_SIZE}/eps{i_episode}") if i_episode % 50 == 0 else gym.make("LunarLander-v2",
19         render_mode="rgb_array")
20     state, info = env.reset()
21     done = False
22     cr = 0
23     while not done:
24         action = agent.take_action(state, eps)
25         next_state, reward, done, truncated, info = env.step(action)
26         agent.experience_replay.store_trans(state, action, next_state, reward, done or truncated)
27         agent.update_params()
28         state = next_state
29         cr += reward
30
31     eps = max(eps * eps_decay_rate, eps_end)
32     crs[i_episode - 1] = cr
33     crs_recent.append(cr)
34     if i_episode % 50 == 0:
35         agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
36
37     print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps:.2f}
38         ', end="")
39     if i_episode % 25 == 0:
40         print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps
41             :.2f}')
42
43     fig = plt.figure()
44     plt.plot(np.arange(len(crs)), crs)
45     plt.ylabel('Reward')
46     plt.xlabel('Episode')
47     plt.title(f"DQN_batch{BATCH_SIZE}")
48     plt.savefig(f"DQN_batch{BATCH_SIZE}.pdf")
49     plt.show()
50     os.system("zip -r DQN.zip DQN/")
```



این بخش از کد، حلقه آموزش را برای اندازه دسته ۳۲ اجرا می کند. در هر اپیزود، محیط reset می شود و عامل با استفاده از سیاست e-greedy اقدام می کند. تجربه های عامل در ExperienceReplay ذخیره می شوند و پارامترهای شبکه به روزرسانی می شوند. در هر ۵۰ اپیزود، مدل ذخیره می شود و در پایان، نتایج آموزش به صورت نمودار نمایش داده می شوند.

```
1 env = gym.make('LunarLander-v2', render_mode="rgb_array")
2 state_size = env.observation_space.shape[0]
3 action_size = env.action_space.n
4 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

در این بخش، محیط LunarLander-v2 تعریف می شود. اندازه های حالت ها و اقدامات از محیط استخراج می شوند.

```
1 n_episodes = 250
2 eps = 1.0
3 eps_decay_rate = 0.97
4 eps_end = 0.01
5 BATCH_SIZE = 32
```

در این بخش، پارامترهای آموزشی تنظیم می شوند. تعداد اپیزودها، مقدار اولیه epsilon، نرخ کاهش epsilon، مقدار نهایی epsilon و اندازه دسته (Batch Size) تعیین می شوند.

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

عامل DQN تعریف می شود و متغیرهایی برای ذخیره پاداش های کسب شده در هر اپیزود و پاداش های اخیر ایجاد می شوند.

```
1 for i_episode in range(1, n_episodes+1):
2     env = RecordVideo(gym.make("LunarLander-v2", render_mode="rgb_array"), f"./DQN/batch{
3         BATCH_SIZE}/eps{i_episode}") if i_episode % 50 == 0 else gym.make("LunarLander-v2",
4         render_mode="rgb_array")
5     state, info = env.reset()
6     done = False
7     cr = 0
8     while not done:
9         action = agent.take_action(state, eps)
10        next_state, reward, done, truncated, info = env.step(action)
11        agent.experience_replay.store_trans(state, action, next_state, reward, done or truncated)
12        agent.update_params()
13        state = next_state
14        cr += reward
15
16    eps = max(eps * eps_decay_rate, eps_end)
17    crs[i_episode - 1] = cr
18    crs_recent.append(cr)
19    if i_episode % 50 == 0:
20        agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
21
22    print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps:.2f}',
23          , end="")
```



```

21 if i_episode % 25 == 0:
22     print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps:.2f}')

```

در این بخش، حلقه اصلی آموزش اجرا می‌شود:

- در ابتدای هر اپیزود، محیط reset می‌شود و حالت اولیه تنظیم می‌شود.
- در هر گام زمانی، عامل یک اقدام را با استفاده از سیاست e-greedy انتخاب می‌کند.
- اقدام انتخاب شده اجرا می‌شود و حالت جدید، پاداش و وضعیت اتمام (done) از محیط بازگردانده می‌شوند.
- انتقال‌ها در حافظه Experience Replay ذخیره می‌شوند.
- پارامترهای شبکه عصبی به‌روزرسانی می‌شوند.
- حالت فعلی به حالت جدید به‌روزرسانی می‌شود و پاداش کسب شده در این گام به پاداش کلی اپیزود اضافه می‌شود.

```

1 eps = max(eps * eps_decay_rate, eps_end)
2 crs[i_episode - 1] = cr
3 crs_recent.append(cr)

```

در پایان هر اپیزود، مقدار epsilon کاهش می‌یابد و پاداش کسب شده در این اپیزود ذخیره می‌شود.

```

1 if i_episode % 50 == 0:
2     agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
3
4 fig = plt.figure()
5 plt.plot(np.arange(len(crs)), crs)
6 plt.ylabel('Reward')
7 plt.xlabel('Episode')
8 plt.title(f"DQN_batch{BATCH_SIZE}")
9 plt.savefig(f"DQN_batch{BATCH_SIZE}.pdf")
10 plt.show()
11 os.system("zip -r DQN.zip DQN/")

```

مدل یادگیری شده در فواصل مشخص (هر ۵۰ اپیزود) ذخیره می‌شود و در پایان، نتایج آموزش به صورت نمودار نمایش داده می‌شوند.

حلقه آموزش برای DDQN

```

1 env = gym.make('LunarLander-v2', render_mode="rgb_array")
2 state_size = env.observation_space.shape[0]
3 action_size = env.action_space.n
4 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
5
6 n_episodes = 250
7 eps = 1.0

```



```
8 eps_decay_rate = 0.97
9 eps_end = 0.01
10 BATCH_SIZE = 64
11
12 agent = DDQNAgent(state_size, action_size, batch_size=BATCH_SIZE)
13 crs = np.zeros(n_episodes)
14 crs_recent = deque(maxlen=25)
15
16 for i_episode in range(1, n_episodes+1):
17     env = RecordVideo(gym.make("LunarLander-v2", render_mode="rgb_array"), f"./DDQN/batch{
18         BATCH_SIZE}/eps{i_episode}") if i_episode % 50 == 0 else gym.make("LunarLander-v2",
19         render_mode="rgb_array")
20     state, info = env.reset()
21     done = False
22     cr = 0
23     action_count = 0
24     while not done:
25         action = agent.take_action(state, eps)
26         next_state, reward, done, truncated, info = env.step(action)
27         agent.experience_buffer.store_trans(state, action, next_state, reward, done or truncated)
28         agent.update_params()
29         state = next_state
30         cr += reward
31         action_count += 1
32         if action_count % 5 == 0:
33             agent.update_target_network()
34
35     eps = max(eps * eps_decay_rate, eps_end)
36     crs[i_episode - 1] = cr
37     crs_recent.append(cr)
38     if i_episode % 50 == 0:
39         agent.save(f"q_net_batch{BATCH_SIZE}_eps{i_episode}.pt")
40
41     print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps:.2f}
42         ', end="")
43     if i_episode % 25 == 0:
44         print(f'\rEpisode {i_episode}\tAverage Reward: {np.mean(crs_recent):.2f}\tEpsilon: {eps
45             :.2f}')
46
47 fig = plt.figure()
48 plt.plot(np.arange(len(crs)), crs)
49 plt.ylabel('Reward')
50 plt.xlabel('Episode')
51 plt.title(f"DDQN_batch{BATCH_SIZE}")
52 plt.savefig(f"DDQN_batch{BATCH_SIZE}.pdf")
```

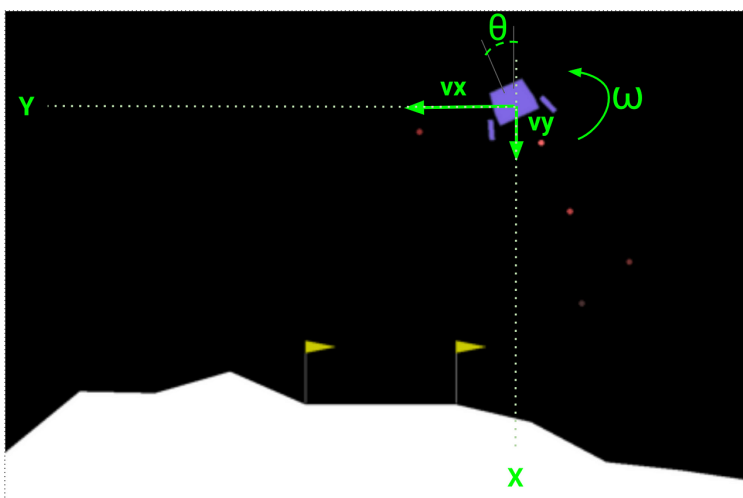


```
49 plt.show()
50 os.system("zip -r DDQN64.zip DDQN/")
```

این بخش از کد، حلقه آموزش را برای DDQN با اندازه دسته ۶۴ اجرا می کند. در هر اپیزود، محیط بازنشانی می شود و عامل با استفاده از سیاست e-greedy اقدام می کند. تجربه های عامل در ExperienceReplay ذخیره می شوند و پارامترهای شبکه به روزرسانی می شوند. در هر ۵۰ اپیزود، مدل ذخیره می شود و در پایان، نتایج آموزش به صورت نمودار نمایش داده می شوند.

۱.۲

این محیط یکی از محیط های Box2D در Gymnasium است که هدف آن شبیه سازی فرود یک فرودگر روی سطح ماه است. این محیط شامل یک بردار مشاهده (observation) هشت بعدی است که مختصات x و y ، سرعت های خطی در x و y زاویه، سرعت زاویه ای، و دو مقدار بولی که نشان می دهند آیا هر پا با زمین در تماس است یا نه، را شامل می شود.



شکل ۱: محیط Lunar lander

این محیط یک محیط پیاده سازی شده توسط Box2D است که به شما امکان می دهد یک فرودگر را به طور مستقل کنترل کنید و آن را با موفقیت روی سطح ماه فرود بیاورید. فضای عمل (action) شامل چهار عمل گسسته است:

۱. عدم انجام کاری
۲. فعال کردن موتور سمت چپ
۳. فعال کردن موتور اصلی
۴. فعال کردن موتور سمت راست

فضای مشاهده (observation) شامل یک بردار هشت بعدی است که شامل موارد زیر است:

- مختصات فرودگر در x و y
- سرعت های خطی در x و y



- زاویه

- سرعت زاویه‌ای

- دو مقدار بولی که نشان می‌دهند هر پا فرودگر با زمین در تماس است یا نه

پاداش‌ها به گونه‌ای طراحی شده‌اند که فرود با موفقیت و استفاده بهینه از سوخت را تشویق کنند. جریمه‌ها برای برخورد با سطح با سرعت بالا، سقوط یا مصرف زیاد سوخت اعمال می‌شوند.

پاداش‌ها

- پس از هر گام، (step) پاداشی اعطا می‌شود. مجموع پاداش یک قسمت جمع تمام پاداش‌ها برای همه گام‌های آن قسمت است.
- برای هر گام، پاداش با نزدیکتر شدن/دور شدن فرودگر به/از پد فرود افزایش/کاهش می‌یابد.
- پاداش با کندتر/سریعتر شدن حرکت فرودگر افزایش/کاهش می‌یابد.
- پاداش با کج شدن فرودگر (زاویه غیر افقی) کاهش می‌یابد.
- برای هر پا که با زمین تماس دارد، ۱۰ امتیاز پاداش افزایش می‌یابد.
- با هر فریم که یک موتور جانبی روشن است، 0.03 امتیاز کاهش می‌یابد.
- با هر فریم که موتور اصلی روشن است، 0.3 امتیاز کاهش می‌یابد.
- هر قسمت یک پاداش اضافی به ترتیب +100 یا -100 امتیاز برای سقوط یا فرود ایمن دریافت می‌کند.
- یک قسمت به عنوان یک راه حل در نظر گرفته می‌شود اگر حداقل ۲۰۰ امتیاز کسب کند.
- شروع و پایان هر قسمت (episode) به صورت زیر تعریف شده است.

شروع حالت

فرودگر از بالای مرکز صفحه با نیروی اولیه تصادفی به مرکز جرم خود شروع می‌کند.

خاتمه قسمت

قسمت به پایان می‌رسد اگر:

- فرودگر سقوط کند (بدنه فرودگر با ماه تماس بگیرد).
- فرودگر از محدوده صفحه خارج شود (مختصات x بزرگتر از ۱ شود).
- فرودگر بیدار نباشد. از مستندات، Box2D بدنه‌ای که بیدار نیست بدنه‌ای است که حرکت نمی‌کند و با هیچ بدنه دیگری برخورد نمی‌کند.



برای استفاده از محیط پیوسته، باید آرگومان `continuous=True` را به شکل زیر مشخص کنید:

```
1 import gymnasium as gym
2 env = gym.make(
3     "LunarLander-v2",
4     continuous: bool = False,
5     gravity: float = -10.0,
6     enable_wind: bool = False,
7     wind_power: float = 15.0,
8     turbulence_power: float = 1.5,
9 )
```

اگر `continuous=True` تنظیم شود، اعمال پیوسته (مربوط به قدرت موتور) استفاده می‌شود و فضای عمل $(-1, +1, (2,))$ Box(-1, +1, (2,)) خواهد بود. مختصات اول عمل، قدرت موتور اصلی و مختصات دوم، قدرت موتورهای جانبی را تعیین می‌کند. اگر آرایه `np.array([main, lateral])` باشد، موتور اصلی کاملاً خاموش خواهد شد اگر `main < 0` و قدرت موتور از 50% تا 100% برای `0 <= main <= 1` تنظیم می‌شود. به همین ترتیب، اگر $-0.5 < lateral < 0.5$ باشد، موتورهای جانبی روشن نمی‌شوند. اگر $lateral < -0.5$ باشد، موتور چپ و اگر $lateral > 0.5$ باشد، موتور راست روشن می‌شود. باز هم، قدرت موتور از 50% تا 100% بین -1 و -0.5 (و 0.5 و 1) تنظیم می‌شود.

`gravity` مقدار ثابت گرانش را تعیین می‌کند که بین 0 و -12 تنظیم می‌شود.

اگر `enable_wind=True` تنظیم شود، اثرات باد به فرودگر اعمال می‌شود. باد با استفاده از تابع

$$\tanh(\sin(2k(t + C)) + \sin(\pi k(t + C)))$$

تولید می‌شود. `k` به 0.01 تنظیم شده و `C` به طور تصادفی بین 9999- و 9999+ نمونه برداری می‌شود.

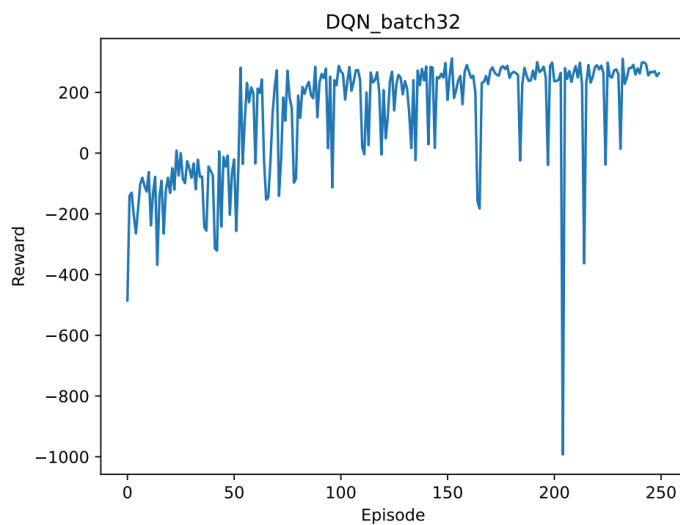
`wind_power` حداکثر اندازه باد خطی اعمال شده به وسیله نقلیه را تعیین می‌کند. مقدار پیشنهادی برای `wind_power` بین 0.0 و 20.0 است. `turbulence_power` حداکثر اندازه باد چرخشی اعمال شده به وسیله نقلیه را تعیین می‌کند. مقدار پیشنهادی برای `turbulence_power` بین 0.0 و 2.0 است.

تاریخچه نسخه‌ها

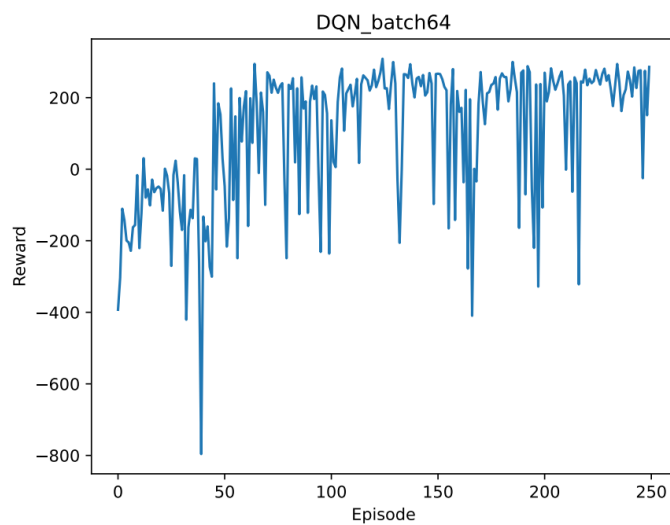
- نسخه v2: شمارش انرژی مصرفی و در نسخه 0.24، اضافه شدن توربولانس با پارامترهای `wind power` و `turbulence_power`.
- نسخه v1: اضافه شدن تماس پاها با زمین در بردار حالت؛ تماس با زمین 10 امتیاز پاداش و از دست دادن تماس 10 امتیاز جریمه؛ پاداش به 200 بازنرمالیزه شده؛ فشار اولیه تصادفی سخت‌تر.
- نسخه v0: نسخه اولیه.

۲.۲

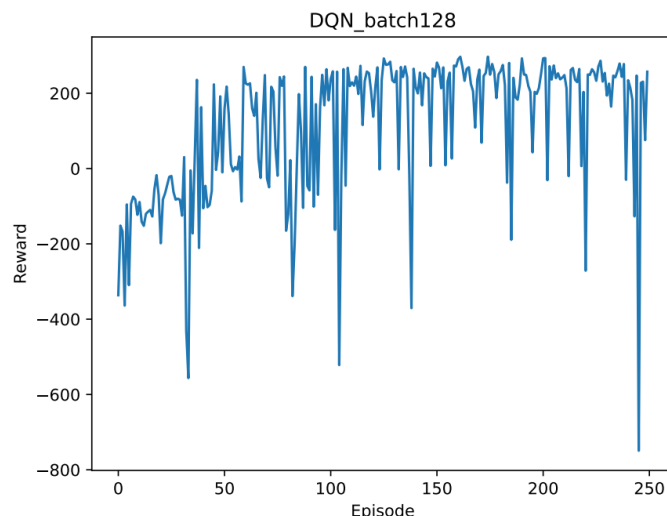
با پیاده سازی این کد برای `batch size` های متفاوت نمودار `reward` بر حسب تعداد `episode` به صورت زیر است.



شکل ۲: نمودار reward برای batch size=32



شکل ۳: نمودار reward برای batch size=64



شکل ۴: نمودار reward برای batch size=128

به طور کلی اندازه دسته در شبکه‌های عصبی عمیق Q (DQN) می‌تواند به طور قابل توجهی بر عملکرد و پایداری فرآیند آموزش تأثیر بگذارد. اندازه‌های دسته کوچک‌تر می‌توانند منجر به به‌روزرسانی‌های پایدارتر در یادگیری شوند زیرا واریانس بیشتری در تخمین گرادیان‌ها ایجاد می‌کنند که می‌تواند به خروج از کمینه‌های محلی کمک کند. با این حال، این واریانس افزایش یافته نیز می‌تواند منجر به ناپایداری در یادگیری شود زیرا به‌روزرسانی‌ها می‌توانند نویزی باشند. در مقابل، دسته‌های بزرگ‌تر تخمین‌های گرادیانی پایدارتر فراهم می‌کنند زیرا میانگین‌گیری بیشتری بر روی نمونه‌ها انجام می‌دهند که منجر به به‌روزرسانی‌های روان‌تر می‌شود. با این حال، ممکن است سرعت خروج از کمینه‌های محلی را کند کنند و به حافظه بیشتری نیاز داشته باشند. اندازه‌های دسته کوچک ممکن است در ابتدا سریع‌تر همگرا شوند به دلیل به‌روزرسانی‌های مکرر، اما به دلیل به‌روزرسانی‌های نویزی همگرایی ممکن است بهینه نباشد. اندازه‌های دسته بزرگ تمایل دارند در ابتدا به آرامی همگرا شوند اما در بلندمدت می‌توانند به یک سیاست بهینه‌تر برسند به دلیل به‌روزرسانی‌های پایدار و دقیق‌تر. دسته‌های کوچک می‌توانند به مدل کمک کنند تا بهتر تعمیم یابد و از بیش‌برازش جلوگیری کند. در مقابل، دسته‌های بزرگ ممکن است منجر به بیش‌برازش شوند زیرا مدل می‌تواند بیش از حد به نمونه‌های خاص در دسته بزرگ وابسته شود. اندازه‌های دسته کوچک می‌توانند از نظر محاسباتی ناکارآمد باشند به دلیل بهینه نبودن استفاده از سخت‌افزار (مثل GPUها). هر گام به‌روزرسانی ممکن است شامل داده کمتری باشد که منجر به کم‌استفاده شدن از توان پردازشی شود. در مقابل، دسته‌های بزرگ بر روی سخت‌افزاری مثل GPUها به دلیل استفاده بهتر از قابلیت‌های پردازش موازی کارآمدترند، اما ممکن است به حافظه و منابع محاسباتی بیشتری نیاز داشته باشند. محدودیت‌های حافظه یکی از ملاحظات عملی مهم است؛ اندازه‌های دسته بزرگ‌تر به حافظه بیشتری نیاز دارند. دسته‌های کوچک‌تر ممکن است در ابتدا به اکتشاف بهتر کمک کنند به دلیل به‌روزرسانی‌های نویزی، در حالی که دسته‌های بزرگ‌تر ممکن است به بهره‌برداری از سیاست‌های یادگرفته شده کمک کنند به دلیل به‌روزرسانی‌های پایدار. اندازه دسته بهینه اغلب نیاز به تعیین تجربی دارد و می‌تواند به مشکل خاص و محیط بستگی داشته باشد.

با توجه به نمودارها و توضیحات فوق مشاهده می‌شود که با دسته ۳۲ نتایج بهتری داریم چراکه تعداد قسمت‌های ما زیاد نیست و پایدار شدن دسته‌های بزرگ‌تر نیاز به زمان و قست‌های بیشتری دارد.

حال به سراغ ارزیابی این سه حالت به کمک متریک regret می‌رویم. ابتدا تعریفی از این معیار را با هم مرور کنیم. این معیار برای اندازه‌گیری اوضاع در شرایط تصمیم‌گیری و در فضای یادگیری تقویتی به کار می‌رود. این معیار تفاوت بین پاداشی که توسط ایجنت دریافت شده و پاداشی که دریافت می‌شد، اگر ایجنت بهترین تصمیم را در هر اپیزود می‌گرفت، ارزیابی می‌کند. پس می‌توان آن را به



صورت رابطه ۱ فرموله کرد.

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

در رابطه ۱، μ^* پاداشی است که اگر ایجنت بهترین اکشن را انتخاب می‌کرد دریافت می‌کرد. در طرف دیگر μ_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 اپیزودها حساب کنیم، می‌توانیم regret تجمعی برای ایجنت را به دست آوریم. در واقع می‌توان گفت regret هزینه انتخاب نکردن بهترین تصمیم در هر اپیزود را به ما نشان می‌دهد. مطابق نتایج، حالت $\text{batch size} = 32$ در ابتدا regret متوسطی دارد، سپس این معیار به صورت واضحی کم می‌شود چون پاداش‌ها در حال افزایش هستند (حدود اپیزود ۵۰) و در نهایت در اپیزودهای بالاتر این معیار به کمترین حالت خود می‌رسد از آنجایی که پاداش‌ها در حال افزایش هستند. برای حالت $\text{batch size} = 128$ ما کمترین میزان regret را در فازهای مختلف شاهد هستیم چون بیشترین پاداش‌ها را برای این حالت دریافت می‌کنیم و می‌توان این‌طور توصیف کرد که در این حالت ایجنت به بهترین عملکرد خود نزدیک است که این مورد با توضیحات قبلی در رابطه با پایداری در دسته‌های بزرگتر و عملکرد مثبت در طولانی مدت همخوانی دارد.

۳.۲

تفاوت اصلی بین Deep Q-Network (DQN) و Double Deep Q-Network (DDQN) در روش به‌روزرسانی Q-values است. در DQN، شبکه عصبی برای تقریب Q-values استفاده می‌شود و برای به‌روزرسانی Q-values از شبکه فعلی استفاده می‌شود. این امر می‌تواند منجر به بیش‌برازش شود زیرا شبکه تمایل دارد به صورت ناپایدار مقادیر Q بزرگ‌تری را برای اقدامات خاص اختصاص دهد. در مقابل، DDQN از دو شبکه عصبی جداگانه استفاده می‌کند: یک شبکه برای انتخاب اقدامات (online network) و یک شبکه برای ارزیابی مقادیر Q (target network). این روش به کاهش بیش‌برازش کمک می‌کند زیرا اقدامات بر اساس شبکه آنلاین انتخاب می‌شوند اما مقادیر Q بر اساس شبکه هدف به‌روزرسانی می‌شوند. این رویکرد منجر به تخمین‌های پایدارتر و دقیق‌تر از مقادیر Q می‌شود و از برخی از مشکلات ناپایداری که در DQN مشاهده می‌شود جلوگیری می‌کند. به عبارت دیگر، DDQN با استفاده از دو شبکه جداگانه برای انتخاب و ارزیابی اقدامات، اطمینان حاصل می‌کند که مقادیر Q به صورت متوازن‌تری به‌روزرسانی می‌شوند و این منجر به بهبود عملکرد و پایداری در فرآیند یادگیری تقویتی می‌شود. این تفاوت اساسی بین DQN و DDQN را به گزینه‌ای مطلوب‌تر برای بسیاری از کاربردهای یادگیری تقویتی تبدیل کرده است، به ویژه در محیط‌هایی که پیچیدگی و ناپایداری بالایی دارند. حال این تفاوت را در کد بررسی می‌کنیم.

- در DQN، تنها یک شبکه عصبی (value_net) برای ارزیابی اقدامات استفاده می‌شود.
- در DDQN، دو شبکه عصبی مجزا وجود دارد: شبکه ارزش (value_net) و شبکه هدف (target_value_net). شبکه هدف به صورت دوره‌ای به‌روزرسانی می‌شود تا از بیش‌برازش جلوگیری شود.

```
1 agent.update_target_network()
```

```
2
```

این خط کد در DDQN استفاده می‌شود تا شبکه هدف به‌روزرسانی شود. در DQN، چنین به‌روزرسانی وجود ندارد.



- در هر دو مدل، (Experience Replay) برای ذخیره انتقال‌ها استفاده می‌شود. با این حال، در DDQN این حافظه با نام `experience_buffer` شناخته می‌شود.

```
1 agent.experience_buffer.store_trans(state, action, next_state, reward, done or truncated)
2
```

این خط کد نشان‌دهنده استفاده از Experience Replay در DDQN است. در DQN، این حافظه با نام `experience_replay` شناخته می‌شود.

- در DQN، به‌روزرسانی پارامترها تنها با استفاده از شبکه ارزش انجام می‌شود.
- در DDQN، به‌روزرسانی پارامترها با استفاده از شبکه ارزش و شبکه هدف انجام می‌شود تا از مشکلات ناشی از تخمین‌های مغرضانه جلوگیری شود.

```
1 q_targets_next = self.target_value_net(next_state_batch).detach().max(1)[0].unsqueeze(1)
2
```

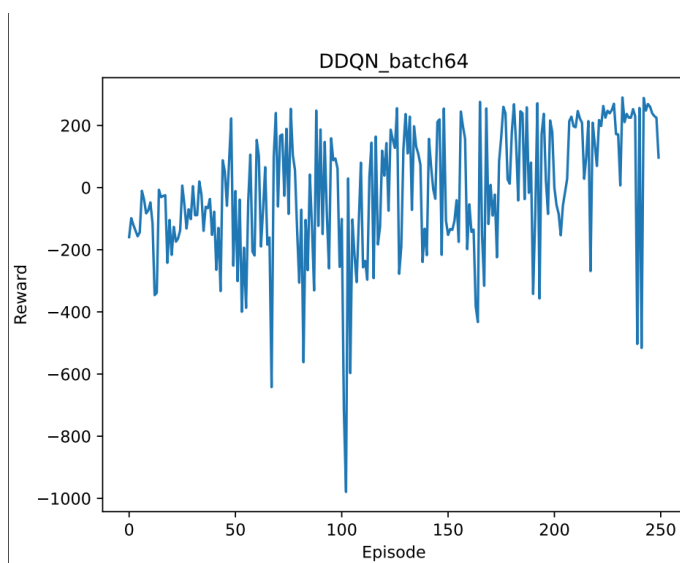
این خط کد نشان‌دهنده استفاده از شبکه هدف در DDQN برای به‌روزرسانی پارامترها است. در DQN، این به‌روزرسانی تنها با استفاده از شبکه ارزش انجام می‌شود.

- در DDQN، شبکه هدف هر چند گام به‌روزرسانی می‌شود (در اینجا هر ۵ اقدام).

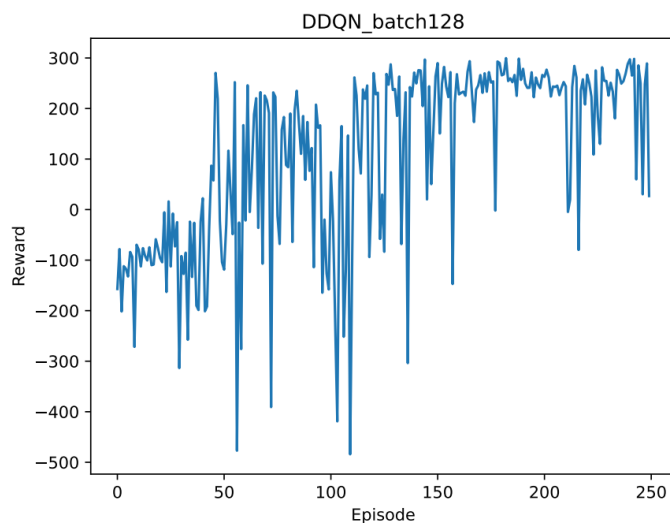
```
1 if action_count % 5 == 0:
2     agent.update_target_network()
3
```

این کد در DDQN استفاده می‌شود تا شبکه هدف هر ۵ اقدام به‌روزرسانی شود. در DQN، چنین فرکانسی برای به‌روزرسانی وجود ندارد.

نتایج پیاده‌سازی با استفاده از DDQN به صورت زیر است.

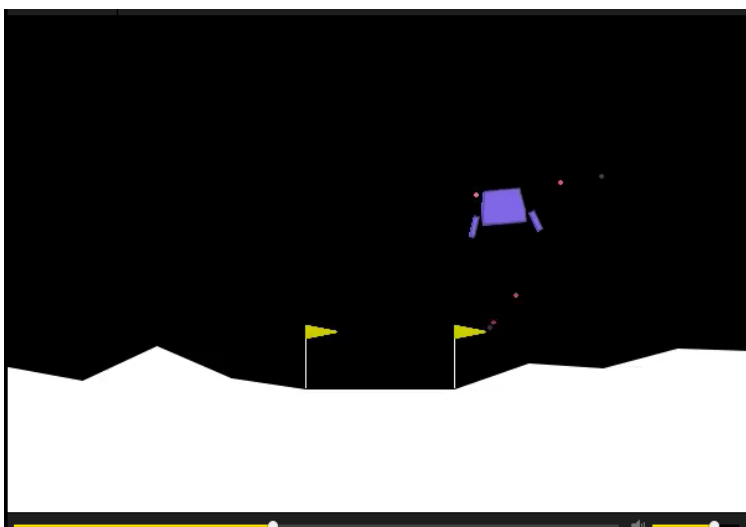


شکل ۵: نمودار reward برای batch size=64



شکل ۶: نمودار reward برای batch size=128

مشاهده می شود که عملکرد در دسته با اندازه ۱۲۸ بهتر است ولی به طور کلی ناپایداری در آن مشاهده می شود. در مقایسه با حالت قبلی، سریعتر همگرا می شود و همچنین از نظر میزان پایداری در پاداش مقداری بهتر عمل می کند هرچند که هنوز نوسان زیادی در آن مشاهده می شود که دلیل آن می تواند پایین بودن تعداد episode ها باشد.



شکل ۷: نمونه خروجی ویدیویی

عملکرد ویدیویی این پیاده‌سازی در تعداد batch size های متفاوت در لینک‌های زیر بارگذاری شده است.

[لینک Github](#)

[لینک Colab](#)