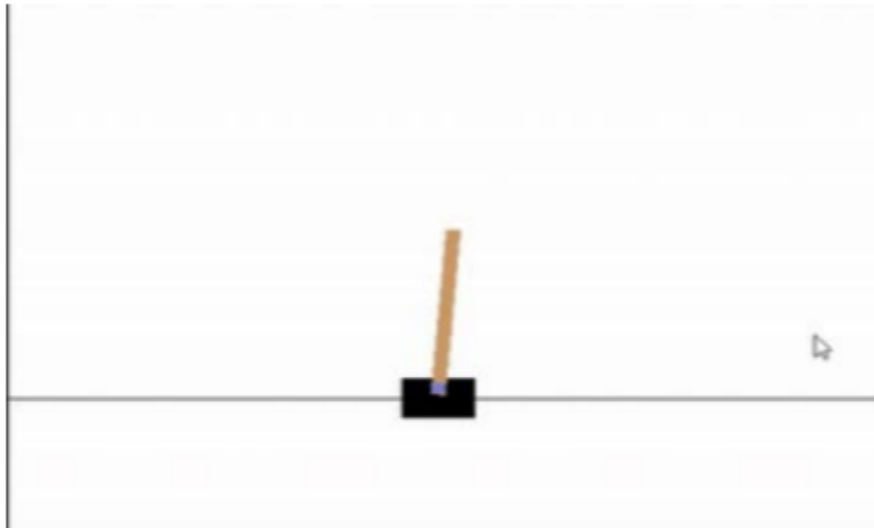


# 1. 作业说明

---

实现倒立摆（CartPole）小游戏，游戏里面有一个小车，上有竖着一根杆子，每次重置后的初始状态会有所不同。小车需要左右移动来保持杆子竖直，为了保证游戏继续进行需要满足以下两个条件：

1. 杆子倾斜的角度 $\theta$ 不能大于 $15^\circ$
2. 小车移动的位置 $x$ 需保持在一定范围（中间到两边各2.4个单位长度）



## 2. 评分标准

---

本次RL小实验设计主要为了学习状态，动作，价值，回报的设计，所以侧重于通过完成报告阐述对该算法了解，因此各分数设计为：

- PDF报告描述该强化学习场景，简要叙述state, action, value, reward设计过程（3）
- 给出CartPole算法伪代码（1）
- 代码、运行结果和截图（1）

## 3. 设计过程

---

### 3.1. 环境描述

---

使用环境是OpenAi Gym中的CartPole-v1环境，但是根据作业要求做了自定义

#### Episode End

The episode ends if any one of the following occurs:

1. Termination: Pole Angle is greater than  $\pm 12^\circ$
2. Termination: Cart Position is greater than  $\pm 2.4$  (center of the cart reaches the edge of the display)
3. Truncation: Episode length is greater than 500 (200 for v0)

作业要求的角度是 $15^\circ$ ，因此在MyModel文件中对环境进行集成做了自定义

## 3.2. 强化学习场景描述

这是一个经典的倒立摆问题。在这个场景中，一个杆子通过一个不可控制的关节连接到一个小车上，这个小车可以沿着无摩擦的轨道移动。杆子被放置在小车上并竖直向上。智能体的目标是通过向小车施加左右方向的力，保持杆子竖直不倒。

## 3.3. 详细设计

### State:

Observation Space

这个环境的状态 (state) 由一个长度为 4 的向量构成，包含了以下信息：

1. Cart Position (小车位置)
2. Cart Velocity (小车速度)
3. Pole Angle (杆子角度)
4. Pole Angular Velocity (杆子角速度)

### Action:

Action Space

供选择的动作是离散的，只有两种：0 或 1。这表示对小车施加的力的方向。

- 0: 向左推小车
- 1: 向右推小车

### Value:

Value Function

在这里，价值 (value) 函数是由 DQN (Deep Q-Network) 网络来近似估计的状态动作值函数。它表示每个状态动作对的预期累积奖励（就是采取这个动作从长期来看收益如何）

### Reward:

可以简单理解为每多撑一轮奖励就会加一

在 CartPole 环境中，每一步操作都会获得 +1 的奖励。目标是尽可能地保持杆子竖直，所以只要杆子没有超出倾斜角度或小车位置超出范围，就会得到奖励 +1。而在该版本的环境中，完成任务所需的奖励阈值为 500。

也就是超过500就会被截断，不过这个参数也是可以改的

## 4. 伪代码

```
1  初始化环境
2  注册定制的CartPole环境，修改角度偏差值
3  创建经验回放池
4  定义DQN网络结构
5
6  for 循环进行多个episode:
7      初始化环境并获取初始状态
8      while 循环进行每个time step:
```

```

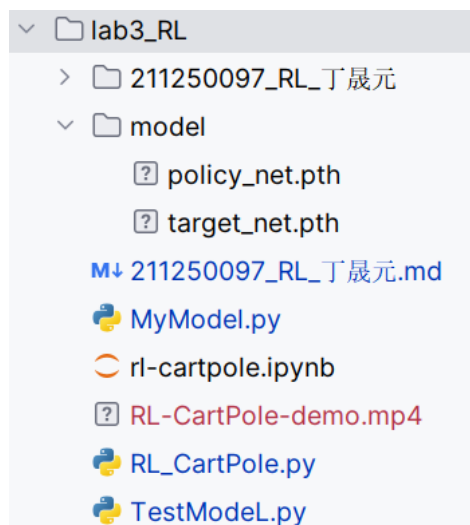
9         通过epsilon-greedy策略选择action（使用DQN网络或随机选择）
10        执行action，获取新状态和奖励
11        存储经验（状态，动作，新状态，奖励）到经验回放池
12        将初始状态更新为新状态
13        从经验回放池中抽取部分经验用于训练DQN网络（通过计算loss并反向传播）
14        如果当前时间步不是最终时间步，则：
15            通过DQN网络估计下一个状态的Q值
16            计算损失（实际的Q值和目标Q值之间的差异）然后优化网络参数
17        更新target network的权重
18        如果达到终止条件，则退出while循环
19    绘制和可视化训练进程
20
21    保存训练好的模型
22    关闭环境

```

## 5. 代码

主要是两个py文件，一个是MyModel.py，一个是RL\_CartPole.py

测试时可以用TestModel.py和model下的pth参数（权重和偏置）文件



### 5.1. MyModel.py

```

1  import math
2  from typing import Optional
3
4  from gymnasium.envs.classic_control import CartPoleEnv
5
6
7  # 创建一个新的 CartPole 环境类，继承自原始的 CartPoleEnv
8  class CustomCartPoleEnv(CartPoleEnv):
9      metadata = {
10          "render_modes": ["human", "rgb_array"],
11          "render_fps": 50,
12      }
13
14      def __init__(self, render_mode: Optional[str] = None):

```

```

15         super(CustomCartPoleEnv, self).__init__(render_mode=render_mode)
16
17         # 修改 theta_threshold_radians 参数为 ±15度
18         self.theta_threshold_radians = 15 * 2 * math.pi / 360
19
20     ## 注册新的 CartPole 环境
21     # gym.envs.register(
22     #     id='CustomCartPole-v1',
23     #     entry_point='xxxxx:CustomCartPoleEnv',
24     #     max_episode_steps=500, # 可以根据需要修改
25     #     reward_threshold=475.0, # 可以根据需要修改
26     # )
27     #
28     ## 创建新的环境实例
29     # custom_env = gym.make('CustomCartPole-v1')
30

```

## 5.2. RL\_CartPole.py

```

1  import os
2
3  import gymnasium as gym
4  import math
5  import random
6  import matplotlib
7  import matplotlib.pyplot as plt
8  from collections import namedtuple, deque
9  from itertools import count
10
11  import torch
12  import torch.nn as nn
13  import torch.optim as optim
14  import torch.nn.functional as F
15
16  # 注册新的 CartPole 环境
17  gym.envs.register(
18      id='CustomCartPole-v1',
19      entry_point='MyModel:CustomCartPoleEnv',
20      max_episode_steps=600, # 可以根据需要修改
21  )
22
23  # 创建新的环境实例
24  env = gym.make('CustomCartPole-v1', render_mode="rgb_array")
25
26  # Transition是一个命名元组，用于表示一个转换，包含四个属性: state, action,
27  # next_state, reward
28  Transition = namedtuple('Transition',
29                          ('state', 'action', 'next_state', 'reward'))
30
31  # ReplayMemory是一个有限长度的存储器，用于存储Agent的经验
32  # 也叫作经验回放池 (Experience Replay Pool)
33  class ReplayMemory(object):
34
35      def __init__(self, capacity):

```

```

36         # deque是一个双端队列，可以从头尾两端添加和删除元素
37         self.memory = deque([], maxlen=capacity)
38
39     def push(self, *args):
40         """Save a transition"""
41         # *args是一个可变参数，可以接受任意多个参数
42         self.memory.append(Transition(*args))
43
44     def sample(self, batch_size):
45         # 选取batch_size个转换
46         return random.sample(self.memory, batch_size)
47
48     def __len__(self):
49         # 自定义实现len函数
50         return len(self.memory)
51
52
53     # Q Network是一个简单的全连接神经网络
54     # 输入是状态，输出是每个动作的Q值
55     class DQN(nn.Module):
56
57         def __init__(self, n_observations_, n_actions_):
58             super(DQN, self).__init__()
59             self.layer1 = nn.Linear(n_observations_, 128)
60             self.layer2 = nn.Linear(128, 128)
61             self.layer3 = nn.Linear(128, n_actions_)
62
63         # forward函数定义了前向传播的运算
64         # 返回值是每个动作的Q值
65         # x可能包含多个样本，每个样本是一个状态，返回的是每个样本对应两个动作的Q值
66         def forward(self, x):
67             x = F.relu(self.layer1(x))
68             x = F.relu(self.layer2(x))
69             return self.layer3(x)
70
71
72     # IPython环境是指在Jupyter Notebook或者Jupyter QtConsole中运行Python代码
73     # 如果是IPython环境，那么is_ipython为True，调用display模块中的display函数显示动画
74     is_ipython = 'inline' in matplotlib.get_backend()
75     if is_ipython:
76         from IPython import display
77
78     # 启用交互模式，可以在动画进行中更新图像
79     plt.ion()
80
81     # 如果有GPU，使用GPU，否则使用CPU
82     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
83
84     # BATCH_SIZE是一个批次的大小，即每次从ReplayMemory中随机选取多少个转换进行训练
85     BATCH_SIZE = 128
86     # GAMMA是折扣因子，用于计算折扣回报
87     # 折扣回报是指未来的奖励的折扣累加和，越靠近1，越重视未来的奖励
88     GAMMA = 0.99
89     # EPS_START是起始的探索率， EPS_END是最终的探索率， EPS_DECAY是探索率的衰减率
90     EPS_START = 0.9
91     EPS_END = 0.05

```

```

92 EPS_DECAY = 1000
93 # TAU是目标网络的更新率，LR是优化器的学习率
94 TAU = 0.005
95 LR = 1e-4
96
97 # n_actions是动作的数量，其实在CartPole的例子中就是2
98 n_actions = env.action_space.n
99 # Get the number of state observations
100 state, info = env.reset()
101 n_observations = len(state)
102
103 # 打印state 和 info
104 print("state:")
105 print(state)
106 print("info:")
107 print(info)
108 print("n_observations:")
109 print(n_observations)
110
111 # policy_net是当前的Q网络，target_net是目标Q网络
112 # 前者用于选择动作，后者用于计算目标Q值
113 policy_net = DQN(n_observations, n_actions).to(device)
114 target_net = DQN(n_observations, n_actions).to(device)
115
116 # 一开始target_net初始化为policy_net的参数
117 target_net.load_state_dict(policy_net.state_dict())
118
119 # 优化器使用AdamW，学习率为LR，amsgrad=True表示使用amsgrad算法
120 optimizer = optim.AdamW(policy_net.parameters(), lr=LR, amsgrad=True)
121 # 经验回放池的容量为10000
122 memory = ReplayMemory(10000)
123
124 # steps_done是一个计数器，用于记录Agent总共与环境交互了多少次
125 steps_done = 0
126
127
128 # select_action函数用于根据当前状态选择动作
129 # 依据epsilon-greedy策略选择动作，即以epsilon的概率随机选择动作，以1-epsilon的概率选择Q值最大的动作
130 def select_action(state_):
131     global steps_done
132
133     sample = random.random()
134     eps_threshold = EPS_END + (EPS_START - EPS_END) * \
135         math.exp(-1. * steps_done / EPS_DECAY)
136     steps_done += 1
137
138     # 如果sample大于eps_threshold，那么以当前Q网络选择动作
139     if sample > eps_threshold:
140         # t.no_grad()表示不需要计算梯度，因为我们不需要对Q网络进行训练，只是用于选择动作
141         with torch.no_grad():
142             # 由于我们的输入是一个批次的状态，所以返回的是每个状态对应的两个动作的Q值
143             # 因此我们要用t.max(1)[1]来获取状态对应的两个动作中Q值最大的那个动作的索引
144             # view(1, 1)是将其变形为一个行向量

```

```

145         return policy_net(state_).max(1)[1].view(1, 1)
146     else:
147         # 否则随机选择动作, 调用env.action_space.sample()函数
148         return torch.tensor([[env.action_space.sample()]], device=device,
dtype=torch.long)
149
150
151 episode_durations = []
152
153
154 def plot_durations(show_result=False):
155     plt.figure(1)
156     durations_t = torch.tensor(episode_durations, dtype=torch.float)
157     if show_result:
158         plt.title('Result')
159     else:
160         plt.clf()
161         plt.title('Training...')
162     plt.xlabel('Episode')
163     plt.ylabel('Duration')
164     plt.plot(durations_t.numpy())
165     # Take 100 episode averages and plot them too
166     if len(durations_t) >= 100:
167         means = durations_t.unfold(0, 100, 1).mean(1).view(-1)
168         means = torch.cat((torch.zeros(99), means))
169         plt.plot(means.numpy())
170
171     plt.pause(0.001) # pause a bit so that plots are updated
172     if is_ipython:
173         if not show_result:
174             display.display(plt.gcf())
175             display.clear_output(wait=True)
176         else:
177             display.display(plt.gcf())
178
179
180 # new
181 # 用来显示CartPole的状态
182 def plot_cartpole_state(screen_):
183     plt.figure(2)
184     plt.clf()
185     plt.imshow(screen_, interpolation='none')
186     plt.title('CartPole State')
187     plt.axis('off')
188     plt.pause(0.001)
189     if is_ipython:
190         display.clear_output(wait=True)
191         display.display(plt.gcf())
192
193
194 # training loop
195 def optimize_model():
196     if len(memory) < BATCH_SIZE:
197         return
198     transitions = memory.sample(BATCH_SIZE)
199

```

```

200     # 将一个批次的转换解压为一个批次的状态，一个批次的动作，一个批次的下一个状态，一个批
    次的奖励
201     batch = Transition(*zip(*transitions))
202
203     # 首先计算一个批次的非终止状态的掩码，即next_state不为None就为True，否则为False
204     non_final_mask = torch.tensor(tuple(map(lambda s: s is not None,
    batch.next_state)), device=device,
205                                     dtype=torch.bool)
206     # 提取了这个批次中所有next_state不为None的转换，然后将其拼接成一个tensor
207     non_final_next_states = torch.cat([s for s in batch.next_state if s is
    not None])
208     state_batch = torch.cat(batch.state)
209     action_batch = torch.cat(batch.action)
210     reward_batch = torch.cat(batch.reward)
211
212     # 将当前状态输入policy_net，得到每个状态对应的两个动作的Q值
213     # 再使用gather函数，将action_batch中的动作对应的Q值提取出来
214     state_action_values = policy_net(state_batch).gather(1, action_batch)
215
216     # 计算下一个状态的Q值，首先初始化为0
217     next_state_values = torch.zeros(BATCH_SIZE, device=device)
218     with torch.no_grad():
219         next_state_values[non_final_mask] =
    target_net(non_final_next_states).max(1)[0]
220     # 预期的值是下一个状态的Q值乘以折扣因子再加上奖励  $r + \gamma * \max_a' Q(s', a')$ 
221     expected_state_action_values = (next_state_values * GAMMA) +
    reward_batch
222
223     # 计算损失，使用smooth_l1_loss函数，即Huber Loss
224     criterion = nn.SmoothL1Loss()
225     loss = criterion(state_action_values,
    expected_state_action_values.unsqueeze(1))
226
227     # 优化模型，首先将梯度置为0，然后进行反向传播，最后进行梯度裁剪
228     optimizer.zero_grad()
229     loss.backward()
230     torch.nn.utils.clip_grad_value_(policy_net.parameters(), 100)
231     # 使用优化器来更新模型参数，使得梯度下降
232     optimizer.step()
233
234
235     if torch.cuda.is_available():
236         num_episodes = 1000
237     else:
238         num_episodes = 600
239
240     for i_episode in range(num_episodes):
241         # 初始化环境并获取初始状态
242         # 按照官方文档的说法可以设置初始状态的生成范围
243         state, info = env.reset()
244         # unsqueeze(0) 将原始状态数据变成了一个形状为 (1, ...) 的张量，其中 1 表示批次大
    小
245         # 便于后续模型输入
246         state = torch.tensor(state, dtype=torch.float32,
    device=device).unsqueeze(0)
247

```



```

248 # 这是一种无限循环的写法，直到任务终止才会退出循环
249 for t in count():
250     # 根据状态选择动作
251     action = select_action(state)
252     # step函数执行动作，返回新的状态、奖励、是否终止、是否被截断等信息
253     print(action)
254     print(type(action))
255     observation, reward, terminated, truncated, _ =
env.step(action.item())
256     reward = torch.tensor([reward], device=device)
257     # done是终止或者被截断的标志
258     done = terminated or truncated
259
260     if terminated:
261         next_state = None
262     else:
263         next_state = torch.tensor(observation, dtype=torch.float32,
device=device).unsqueeze(0)
264
265     # 存储状态转换信息到经验回放缓冲区
266     memory.push(state, action, next_state, reward)
267
268     # 切换到下一个状态
269     state = next_state
270
271     # 执行一次模型优化（policy网络的训练）
272     # 取样一个批次的转换，然后进行模型优化
273     optimize_model()
274
275     # new
276     # 获取当前环境的屏幕状态
277     screen = env.render()
278     plot_cartpole_state(screen)
279
280     # 软更新目标网络的权重
281     # 这里的软更新是指将目标网络的权重向策略网络的权重靠近一点
282     target_net_state_dict = target_net.state_dict()
283     policy_net_state_dict = policy_net.state_dict()
284     for key in policy_net_state_dict:
285         target_net_state_dict[key] = policy_net_state_dict[key] * TAU +
target_net_state_dict[key] * (1 - TAU)
286         target_net.load_state_dict(target_net_state_dict)
287
288     if done:
289         episode_durations.append(t + 1)
290         plot_durations()
291         break
292
293 print('Complete')
294 plot_durations(show_result=True)
295 plt.ioff()
296 plt.show()
297 # 结束后关闭环境
298 env.close()
299
300 # 保存模型

```

```
301 print("Saving model...")
302 # 创建文件夹名为model
303 if not os.path.exists('./model'):
304     os.mkdir('./model')
305 # 保存模型参数
306 torch.save(policy_net.state_dict(), './model/policy_net.pth')
307 torch.save(target_net.state_dict(), './model/target_net.pth')
308 print("Model saved!")
```

## 6. 结果

---

运行TestModel.py可以直接看到动画

也可以看本文件夹下的video



RL-CartPole-demo.mp4

2023/12/8 10:56

MP4 文件

1,085 KB

## 7. 参考链接

---

1. 依旧是官网的demo

[Reinforcement Learning \(DQN\) Tutorial — PyTorch Tutorials 2.1.0+cu121 documentation](#)

2. gym官网文档说明

<https://gymnasium.farama.org/>

[https://github.com/Farama-Foundation/Gymnasium/blob/main/gymnasium/envs/classic\\_control/cartpole.py](https://github.com/Farama-Foundation/Gymnasium/blob/main/gymnasium/envs/classic_control/cartpole.py)

CartPole

[Cart Pole - Gymnasium Documentation \(farama.org\)](#)

