

1、定义算法

教程中提到低于**Q learning**，**DQN**本质上是为了适应更为复杂的环境，并且经过不断的改良迭代，到**T Nature DQN**（即Volodymyr Mnih发表的Nature论文）这里才算是基本完善。**DQN**主要改动的点有三个：

- 使用深度神经网络替代原来的Q表：这个很容易理解原因
- 使用了经验回放（Replay Buffer）：这个好处有很多，一个是使用一堆历史数据去训练，比之前用一次就扔掉好多了，大大提高样本效率。另外一个方面是面试题所提的，减少样本之间的相关性，原则上获取经验跟学习阶段是分开的，原来时序的训练数据有可能是不稳定的，打完之后再学习有损是训练的稳定性。跟深度学习中间划分训练测试集时打乱样本是一个道理。
- 打了折扣：如果神经网络有自循环，每隔若干步才打更多更新的策略网络参数更新给旧网络，这样微也是为了训练的稳定性，避免Q值的累计发散，想一下，如果当前有个transition，这个**Q learning**内经过的，一定是记住！！！样本数据对Q值进行了较差的估计，如果接下来从经验回放中抽取的样本正好是连续几个都这样的，很有可能导致Q值的发散（它的青春小一去不回来了）。再打个比方，我们玩PAC或者有关类游戏，有些人为了赢花很多时间去Save and Load，只要出了招，我不管我就加救之前的存档，假设不允许加载呢，就按**DDQN**方法一样训练过程中会受不了，这时候是不是搞两个档，一个档每帧都存一下，另外一个档打了不错的结果再存，也就是若干个间隔再存一下，到最后用间隔若干步数再存的档一般都比每帧都存的档好表现，当然你也可以再搞更多个档，也就是**DQN**增加多个目标网络，但是对于**DDQN**没有多大必要，多几个网络效果不见得会好很多。

1.1、定义模型

前面说了**DQN**的模型不再是Q表，而是一个深度神经网络，这里我只用了一个三星的金迷魂网络（FCN），这种网络也叫多感知机（MLP），至于怎么用**Torch**写网络这里就不多说明了，以下仅供参考。

```
In [9]: import torch.nn as nn
import torch.nn.functional as F
class MLP(nn.Module):
    def __init__(self, n_states,n_actions,hidden_dim=128):
        """初始化Q网络，为全连接网络"""
        super(MLP, self).__init__()
        self.fc1 = nn.Linear(n_states, hidden_dim) # 输入层
        self.fc2 = nn.Linear(hidden_dim,hidden_dim) # 隐藏层
        self.fc3 = nn.Linear(hidden_dim, n_actions) # 输出层

    def forward(self, x):
        # 对当前网络的输出
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        return self.fc3(x)
```

1.2、定义经验回放

经验回放首先是一个有一定容量的，只有存储一定的transition网络才会更新，否则就返回之前的逐步更新了。另外写经验回放的时候一般都需要两个功能或方法，一个是push，即将一个transition样本按照序放到经验回表中，如果满了就把最开始放进去的样本挤掉，因此如果大家学过数据结构的哈希桶用列表来写，虽然这里不足，另外一个就是sample，很简单就是随机来采样一个或者若干个（具体多少就是batch_size）样本供**DQN**网络更新，功能讲清楚了，大家可以按照自己的想法用代码来实现，参考如下。

```
In [18]: from collections import deque
import random
class ReplayBuffer(object):
    def __init__(self, capacity: int) -> None:
        self.capacity = capacity
        self.buffer = deque(maxlen=self.capacity)#deque(maxlen=capacity): 双端队列的“固定容量”特性很关键 —— 当存储的经验数量超过 capacity 时，会自动删除队头元素（最早存入）的经验，保证缓冲区不会无限膨胀。
    def push(self,transitions):
        """存储transition到经验回表中"""
        self.buffer.append(transitions)#每次新物体与环境交互后，调用push 把这次经验存入缓冲区。
    def sample(self, batch_size: int, sequential: bool) -> False:
        if batch_size > len(self.buffer): # 如果批量大小大于经验回表的容量，那就经验回表的容量
            batch_size = len(self.buffer)
        if sequential: # 顺序采样
            rand = random.randint(0, len(self.buffer) - batch_size)
            batch = [self.buffer[i] for i in range(rand, rand + batch_size)]
            return list(batch)
        else: # 随机采样
            batch = random.sample(self.buffer, batch_size)
            return list(batch)
    def clear(self):
        self.buffer.clear()
    def __len__(self):
        """返回当前存储的量"""
        return len(self.buffer)
```

1.3、真定义算法

到了真定义一点的算法，定义算法就比较麻烦，要先定义一些子模块，再定义好子模块之后我们就可以实现我们的算法核心部分，如下，可以看到，其实去掉子模块的话，**DDQN**Q learning的算法结构没啥区别，当然因为神经网络一般需要**Torch**或者**Tensorflow**来写，因此建议大家先去学一学这些工具，比如“[cat_pytorch_in_20_days](#)”。

这里我们主要分析一下**DQN**的更新过程，也就是update函数，首先我们知道目前所有基于深度神经网络更新方式都是梯度下降，如下：

$$\theta_t \leftarrow \theta_t - \nabla_{\theta} L_t(\theta_t)$$

那么怎么求这个L呢，注意到前面我们讲的**DDQN**Q learning算法的一个主要区别就是使用神经网络替代了Q表，而这个 θ 实际上就是神经网络的参数，通常用**Q** ($s_t, a_t; \theta$)表示。根据强化学习的原理我们需要优化的是对应状态下不同动作的长期价值，然后每次选择价值最大且能启动的就能完成一条完整路径，使用神经网络表示Q表时也是如此，我们将输入的状态数作为神经网络的输入层，动作数作为输出层，这样的神经网络的功能就跟在**Q learning**中的Q表是一样的，只不过具有更强的鲁棒性。

讲完了怎么优化化的是这个参数 θ ，接下来我们从代码层面进一步剖析，稍微了解一点**Torch**知识的同学都知道，上面的公式其实只需要定义一个优化器，然后计算损失之后用优化器迭代即可，如下：

```
optimizer = optim.Adam(Q_net.parameters()), lr=0.01 # 定义优化器，对应的网络是Q_net，学习率为0.01
loss = ... # 计算损失
loss.backward() # 然后优化器先zero_grad(), loss再反向传播，然后优化器step()，这是一个固定的套路
```

当然按照这个设计了解一下深度学习中的梯度下降，并且使用numpy实现，这样就会更加清楚整个梯度下降过程到底是怎么回事，上述只是在同学们了解了梯度下降的具体实现方式的前提下为了方便学习更多其他的知识形成的套路，这好比我们玩一个竞速游戏，如果我们之前从来没有玩过竞速游戏，那么肯定是从看教程开始，每个技能一步一步地学起打好基础，然后再再学习技能连接等等也就是形成固定的套路，但是如果不打基础，直接学习套路可能会是一脸懵逼的状态，尤其是很多高阶玩家会对这些连招套路化名称比如光速Qa和123323等等，一开始我们是很懵懂的，等当我们先打好基础，然后再再学习了很多套路之后会发现这些基础并不能用得上，甚至有的时候可能会忽然忘了这些基础，但其实我们并没有忘记，再回顾一遍也能很快找起来，在这点上我想强调的是基础固然重要，但是不要死搬硬套，除非非学术研究需要，再比如我们小学学完简单加减乘除之后很快就去背九九乘法表，而不会去过多纠结一加一等于几的问题，上大学后也是如此，只是很多时候我们可能看起这个什么问题值得研究，但意识不到自己就是在纠结一加一等于几的问题，这也是我在教众基础学习讨论的过程中在论坛上发现的问题。

回归主题，后面的学生会发现数学公式和代码的对应是有一定的规律的，只要通过多加练习跨越了这个壁垒，那么对于往后我们想要更理论化也会轻松许多。我们目前讲了参数的更新过程，但是最关键的是损失是如何计算的，在**DQN**中损失的计算相对来说比较清晰，如下：

$$L(\theta) = (y_t - Q(s_t, a_t; \theta))^2$$

这里的 y_t 通常称为期望值， $Q(s_t, a_t; \theta)$ 称为实际值，这个损失在深度学习通常被称为均方差损失，也就是mseloss，使用这个损失函数通常遵循数学上的最小二乘法，感兴趣的同学可以了解一下深度学习中的各种损失函数以及各自的使用场景， y_t 在**DQN**中一般表示如下：

$$y_t = \begin{cases} r_t & \text{对于终止状态 } s_{t+1} \\ r_t + \gamma \max_a Q(s_{t+1}, a; \theta') & \text{对于非终止状态 } s_{t+1} \end{cases}$$

该公式的意思是说下一个状态对应的最大Q值作为实际值（因为实际值通常不能直接求得，只能近似），这种做法实际上只是一种近似，可能导致导致设计计算问题，也有一些改善的方法具体可以在后面各种改进的**DQN**算法比如**Double DQN**中看到，在这里我们暂时不要深究为什么要求这个近似似然值，然后注意到这里其实有一一终止状态的判断，因为如果当前状态是终止状态，那么实际上是没有下一个状态的，所以**DQN**干脆直接使用对应的实际值表示Q的实际值。

```
In [11]: import torch
import torch.optim as optim
import math
import numpy as np
class DQN:
    def __init__(self,model,memory,cfg):
        self.n_actions = cfg['n_actions']
        self.device = torch.device(cfg['device'])
        self.gamma = cfg['gamma'] # 奖励的折扣因子
        # e-greedy策略相关参数
        self.sample_count = 0 # 用于Epsilon的衰减计数
        self.epsilon = cfg['epsilon_start']
        self.sample_count = 0
        self.epsilon_start = cfg['epsilon_start']
        self.epsilon_end = cfg['epsilon_end']
        self.epsilon_decay = cfg['epsilon_decay']
        self.batch_size = cfg['batch_size']
        self.policy_net = model.to(self.device) # 策略网络，实时更新，用于选动作，跟贪心
        self.target_net = model.to(self.device) # 目标网络，固定一段时间更新，用于计算目标Q值（稳定策略）
        # 更新参数到目标网络
        for target_param, param in zip(self.target_net.parameters(),self.policy_net.parameters()):
            target_param.data.copy_(param.data)
        self.optimizer = optim.Adam(self.policy_net.parameters()), lr=cfg['lr'] # 优化器
        self.memory = memory # 经验回放
        def sample_action(self, state):
            """ 采样动作 """
            self.sample_count += 1
            # epsilon探索策略
            self.epsilon = self.epsilon_end + (self.epsilon_start - self.epsilon_end) * \
                math.exp(-1 * self.sample_count / self.epsilon_decay)
            if random.random() < self.epsilon:
                with torch.no_grad():
                    # 状态转换:从numpy数组到tensor，加个维度（适配batch格式，如(4,)-(1,4)）
                    state = torch.tensor(state, device=self.device, dtype=torch.float32).unsqueeze(dim=0)
                    q_values = self.policy_net(state)
                    action = q_values.max(1)[1].item() # choose action corresponding to the maximum q value
            else:
                action = random.randrange(self.n_actions)
            return action
        torch.no_grad() # 不计算梯度，训练效果等同于with torch.no_grad().
        def predict_action(self, state):
            """ 预测动作 """
            # 和“利用”逻辑一样，但没有epsilon判断（纯Q值最大的动作，用于训练完成后的测试）
            state = torch.tensor(state, device=self.device, dtype=torch.float32).unsqueeze(dim=0)
            q_values = self.policy_net(state)
            action = q_values.max(1)[1].item() # choose action corresponding to the maximum q value
            return action
        def update(self):
            if len(self.memory) < self.batch_size: # 当经验回表中不满是一个批量时，不更新策略
                return
            # 从经验回表中随机采样一个批量的转换(transition)
            state_batch, action_batch, reward_batch, next_state_batch, done_batch = self.memory.sample(self.batch_size)
            # 将数据转换成tensor
            state_batch = torch.tensor(np.array(state_batch), device=self.device, dtype=torch.float)
            action_batch = torch.tensor(action_batch, device=self.device).unsqueeze(1)
            reward_batch = torch.tensor(reward_batch, device=self.device, dtype=torch.float)
            next_state_batch = torch.tensor(np.array(next_state_batch), device=self.device, dtype=torch.float)
            done_batch = torch.tensor(np.array(done_batch), device=self.device)
            q_values = self.policy_net(state_batch).gather(dim=1, index=action_batch) # 计算当前状态(s_t,a)对应的Q(s_t,a)
            next_q_values = self.policy_net(next_state_batch).max(1)[0].detach() # 计算下一时刻的状态(s_t,a)对应的Q值
            # 计算期望的Q值，对于终止状态，此时done_batch[0]=1，对应expected_q_value等于reward
            expected_q_values = reward_batch + self.gamma * next_q_values # (1-done_batch)
            loss = nn.MSELoss((q_values - expected_q_values).unsqueeze(1)) # 计算均方误差
            # 优化更新模型
            self.optimizer.zero_grad() # 清空上一轮的梯度（避免累积）
            loss.backward() # 反向传播，计算损失对网络参数的梯度
            # clip的止损梯度操作
            for param in self.policy_net.parameters():
                param.grad.data.clamp_(-1, 1)
            self.optimizer.step()
```

$$e = e_{end} + (e_{start} - e_{end}) \times e^{-1 \times sample_count / decay}$$

2、定义训练

```
In [21]: def train(cfg, env, agent):
    """ 训练 """
    print("开始训练！")
    rewards = [] # 记录所有回合的奖励
    steps = []
    for i_ep in range(cfg['train_eps']):
        ep_reward = 0 # 记录一回合内的奖励
        ep_step = 0
        state, _ = env.reset() # 重置环境，返回初始状态
        for _ in range(cfg['ep_max_steps']):
            ep_step += 1
            action = agent.sample_action(state) # 选择动作
            next_state, reward, terminated, truncated, info = env.step(action) # 更新环境，返回transition
            done = terminated or truncated
            agent.memory.push(state, action, reward, next_state, done) # 保存transition
            state = next_state # 更新下一个状态
            agent.update() # 更新策略网络
            ep_reward += reward # 累加奖励
            if done:
                break
            if (i_ep + 1) % cfg['target_update'] == 0: # 策略网络目标网络更新
                agent.target_net.load_state_dict(agent.policy_net.state_dict())
                steps.append(ep_step)
                rewards.append(ep_reward)
            if (i_ep + 1) % 10 == 0:
                print(f"回合: {i_ep+1}/{cfg['train_eps']}，奖励: {ep_reward:.2f}, Epsilon: {agent.epsilon:.3f}")
                print("当前回合:")
                env.close()
                return ['rewards':rewards]

def test(cfg, env, agent):
    """ 测试 """
    rewards = [] # 记录所有回合的奖励
    steps = []
    for i_ep in range(cfg['test_eps']):
        ep_reward = 0 # 记录一回合内的奖励
        ep_step = 0
        state, _ = env.reset() # 重置环境，返回初始状态
        for _ in range(cfg['ep_max_steps']):
            ep_step += 1
            action = agent.predict_action(state) # 选择动作
            next_state, reward, terminated, truncated, info = env.step(action) # 更新环境，返回transition
            done = terminated or truncated
            state = next_state # 更新下一个状态
            ep_reward += reward # 累加奖励
            if done:
                break
            steps.append(ep_step)
            rewards.append(ep_reward)
        print(f"回合: {i_ep+1}/{cfg['test_eps']}，奖励: {ep_reward:.2f}")
        print("当前回合:")
        env.close()
        return ('rewards':rewards)
```

3. 定义环境

```
In [22]: # !pip install gymnasium
import gymnasium as gym
import os
import random
def all_seed(seed = 1):
    """ 所有的seed函数 """
    np.random.seed(seed)
    random.seed(seed)
    torch.manual_seed(seed) # config for CPU
    torch.cuda.manual_seed(seed) # config for GPU
    os.environ['PYTHONHASHSEED'] = str(seed) # config for python scripts
    # config for cudnn
    torch.backends.cudnn.deterministic = True
    torch.backends.cudnn.benchmark = False
    torch.backends.cudnn.enabled = False
def env_agent_config(cfg):
    env = gym.make(cfg['env_name']) # 创建环境
    if cfg['seed'] != 0:
        all_seed(seed=cfg['seed'])
    n_states = env.observation_space.shape[0]
    n_actions = env.action_space.n
    print(f"环境状态维数: {n_states}，动作空间维数: {n_actions}")
    cfg.update({"n_states":n_states,"n_actions":n_actions}) # 更新n_states和n_actions到cfg参数中
    model = MLP(n_states, n_actions, hidden_dim = cfg['hidden_dim']) # 创建模型
    memory = ReplayBuffer(cfg['memory_capacity']) # 创建经验回表
    agent = DQN(model,memory,cfg)
    return env,agent
```

4、设置参数

```
In [23]: import argparse
import matplotlib.pyplot as plt
import seaborn as sns
def get_args():
    """ 超参数 """
    parser = argparse.ArgumentParser(description="hyperparameters")
    parser.add_argument('--env_name',default='CartPole-v0',type=str,help='name of environment')
    parser.add_argument('--algo_name',default='DQN',type=str,help='name of algorithm')
    parser.add_argument('--train_eps',default=200,type=int,help='episodes of training')
    parser.add_argument('--test_eps',default=20,type=int,help='episodes of testing')
    parser.add_argument('--ep_max_steps',default=10000,type=int,help='steps per episode, much larger value can simulate infinite steps')
    parser.add_argument('--gamma',default=0.95,type=float,help='discounted factor')
    parser.add_argument('--epsilon_start',default=0.5,type=float,help='initial value of epsilon')
    parser.add_argument('--epsilon_end',default=0.01,type=float,help='final value of epsilon')
    parser.add_argument('--epsilon_decay',default=500,type=int,help='decay rate of epsilon, the higher value, the slower decay')
    parser.add_argument('--batch_size',default=0.0001,type=float,help='learning rate')
    parser.add_argument('--target_update',default=4,type=int)
    parser.add_argument('--hidden_dim',default=256,type=int)
    parser.add_argument('--device',default='cpu',type=str,help='cpu or cuda')
    parser.add_argument('--seed',default=1,type=int,help='seed')
    args = parser.parse_args()
    # 打印超参数
    print("超参数:")
    print("-" * 100)
    tplt = "-" * 20 + "t" * 20 + "-" * 20
    print(tplt.format("Name", "Value", "Type"))
    for i,v in args.items():
        print(tplt.format(k,v,str(type(v))))
    print("-" * 100)
    return args
def smooth(data, weight=0.9):
    """ 平滑曲线，类似了Tensorboard中的smooth曲线 """
    smoothed = []
    last = data[0]
    for point in data:
        smoothed_val = last * weight + (1 - weight) * point # 计算平滑值
        smoothed.append(smoothed_val)
        last = smoothed_val
    return smoothed
def plot_rewards(rewards, cfg, tag='train'):
    """ 画图 """
    sns.set()
    plt.figure(figsize=(10, 5)) # 创建一个图形实例，方便同时画多个图
    plt.title(f"({tagging curve on {cfg['device']} of {cfg['algo_name']} for {cfg['env_name']})")
    plt.xlabel('episodes')
    plt.plot(rewards, label='rewards')
    plt.plot(smooth(rewards), label='smoothed')
    plt.legend()
    plt.show()
```

5、开始训练

```
In [24]: # 获取参数
cfg = get_args()
# 训练
env, agent = env_agent_config(cfg)
res_dic = train(cfg, env, agent)

# 测试
res_dic = test(cfg, env, agent)
plot_rewards(res_dic['rewards'], cfg, tag='train')

# 画出结果
plot_rewards(res_dic['rewards'], cfg, tag='test')
```

