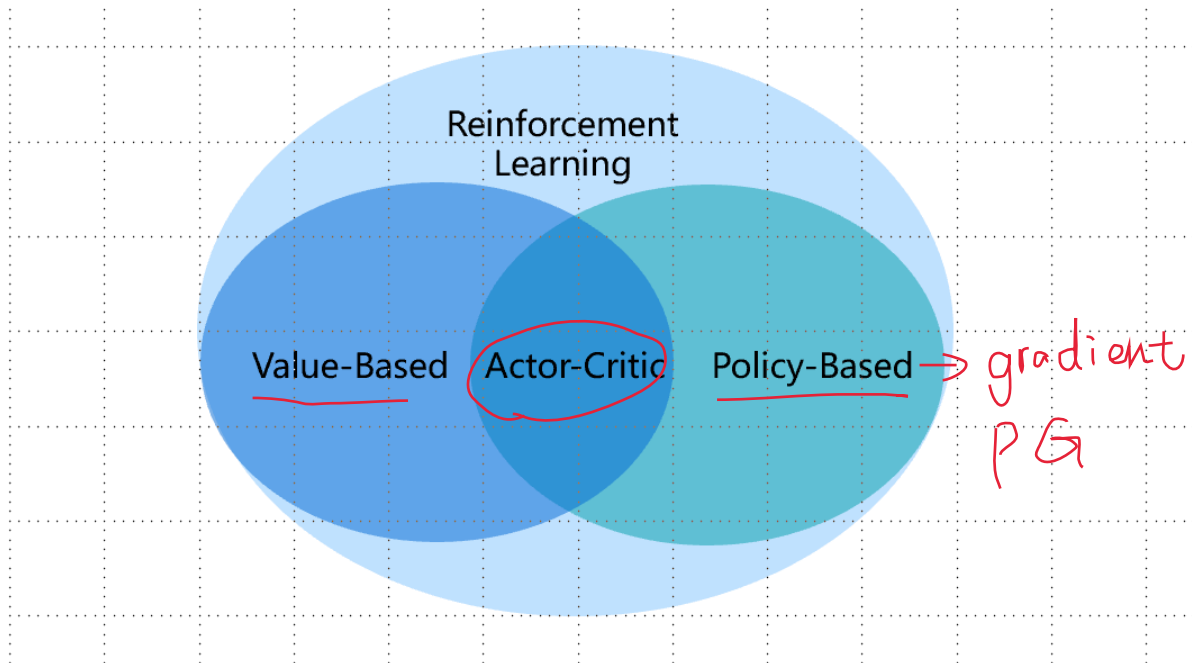


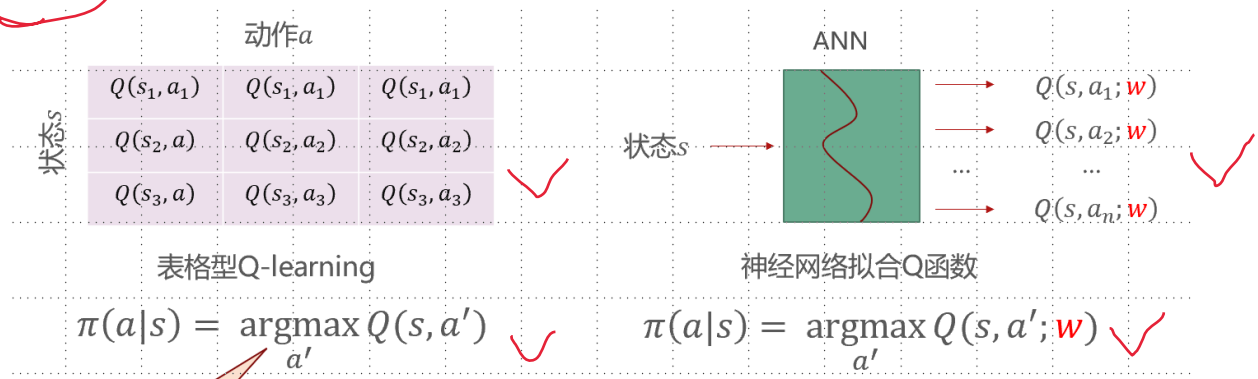
1策略梯度算法——REINFORCE

1 策略梯度算法概述



1.1 基于价值的方法存在的问题

前面介绍的一系列基于价值函数的(value-based)方法，都是估计各个"状态-价值"对的未来收益的期望 $Q(s, a)$ ，然后使用贪婪算法来选择动作。

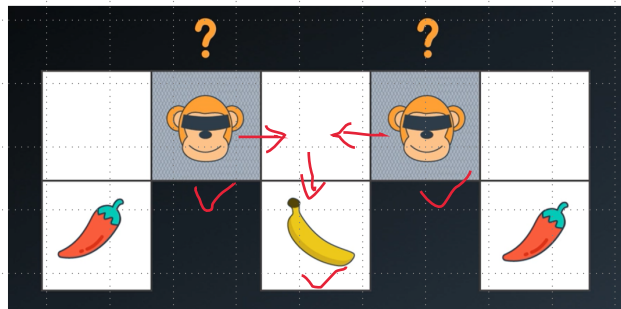
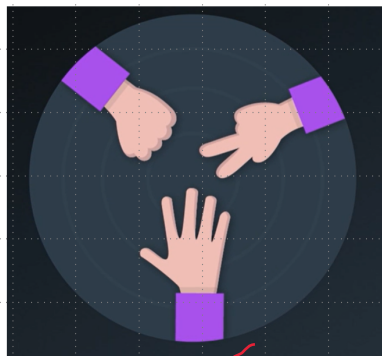


这类方法存在两个比较大的问题：

- 无法产生随机策略

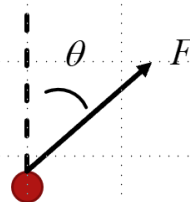
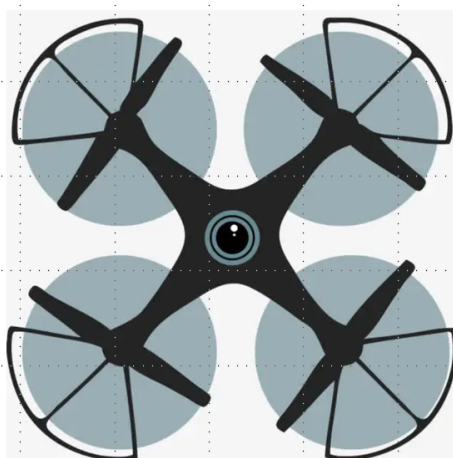
比如玩"石头-剪刀-布"游戏，最优的测率应该是随机策略。如果你有任何喜欢出某一个动作的习惯，都会被对手察觉到，从而对手就会出相应的策略进行压制。基于价值的方法就难以应对这样的情况。

又比如走迷宫，在局部观测条件下，很有可能会遇到两个完全一样的观测，但却有着不同的最优策略。此时同样需要使用随机策略。



- 无法应对连续动作

在很多问题中，智能体的动作在连续空间中变化。比如无人机飞行方向和动力大小的控制问题。



$$\theta \in (-\pi, \pi]$$

$$F \in [0, \infty)$$

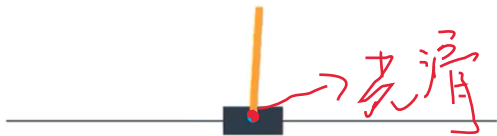
可不可以直接优化智能体的策略呢？当然可以。

智能体的策略函数 π 就是从“状态”到“动作”的映射。即每给定一个状态 s ，策略函数都给出相应的动作 $a = \pi(s)$ 。直接优化智能体的策略是一种更加直观自然的方式，称之为基于策略(policy-based)的方法。

1.2 离散动作和连续动作对应的策略形式

以Cart-Pole为例，小车和杆子的状态为：车的位置、速度、杆子的角度、杆子末端的速度。智能体的动作为：向左、向右。智能体的任务就是不断地根据当前地状态决定向左或者向右施加力，以保持杆子不会倒下。

State

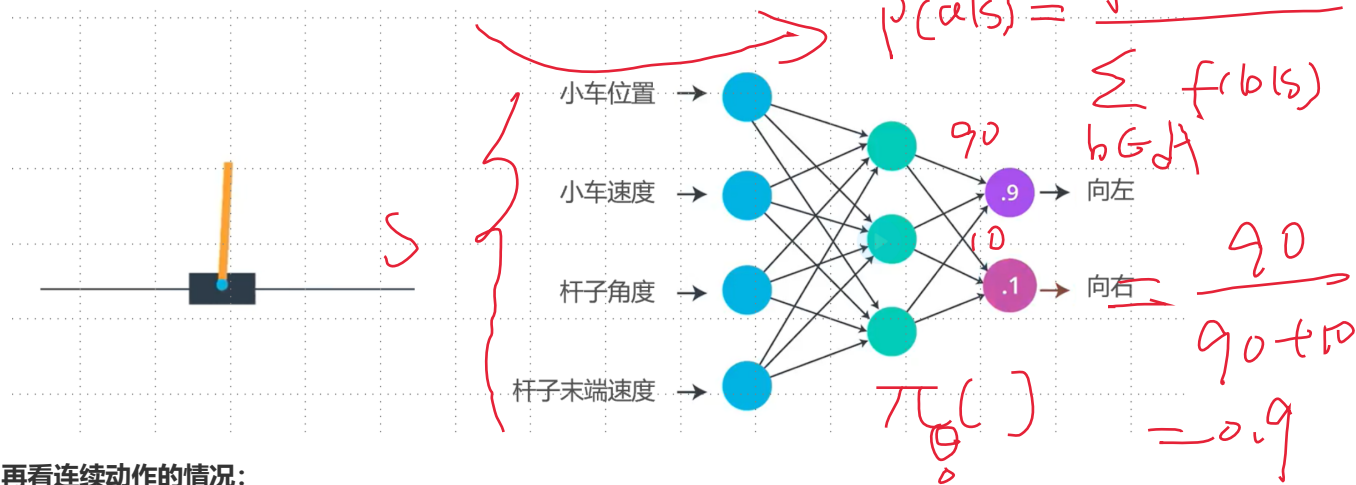


	Min	Max
Cart Position	-2.4	2.4
Cart Velocity	-Inf	Inf
Pole Angle	$\sim -41.8^\circ$	$\sim 41.8^\circ$
Pole Velocity At Tip	-Inf	Inf

首先看离散动作的情况:

如果智能体的动作是离散的, 策略给出每个动作被选择的概率, 即随机策略 $\pi(a|s) = Pr(A = a|S = s)$

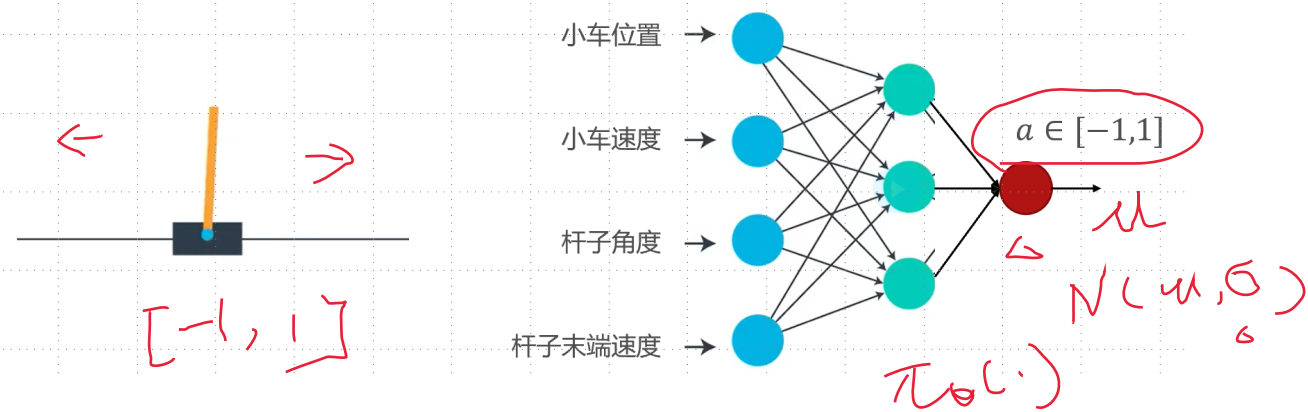
在上面的有图中, 神经网络输出了两个动作的选择概率。为了使得神经网络的输出表示各个动作的概率, 一般对神经网络的输出结果使用Softmax进行处理, 使输出满足概率的特性。



再看连续动作的情况:

如果智能体的动作是连续的, 策略需要给出一个数。这时神经网络的输出不在表示概率, 而是具体的数值。举例来说, 如果Cart-Pole问题中的动作为 $[-1, 1]$ 的一个值, 绝对值表示力的大小, 负表示向左, 正表示向右。如果想限制神经网络的输出永远是合法的, 在神经网络的输出后面加上 \tanh 这样的函数即可。当然, 此时的神经网络最后一层只有一个神经元。

同时, 为了在训练过程中增加对环境的探索, 要对动作引入随机性。具体做的做法是, 神经网络输出动作的均值 μ , 然后构造一个新的正态分布 $N(\mu, \sigma)$, 训练的时候每次都从该分布中采样一个动作。



1.3 策略的参数及其优化

智能体的策略可以用参数化的函数实现，比如神经网络。用 θ 表示策略的参数，则不同的参数唯一确定了智能体与环境交互的策略。为了在交互中获得更好的表现，就必须优化参数 θ （为了方便表述，下面直接称策略为 θ ）。

首先，需要量化当前策略 θ 的表现。

考虑分幕型(episodic)强化学习问题，智能体与环境交互的一条轨迹记为

$$\tau = (s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1}, s_T)$$

这一幕的回报为

$$G(\tau) = r_1 + r_2 + \dots + r_T$$

不同的轨迹对应着不同的回报，那么当前策略 θ 下的期望回报可以记为

$$J(\theta) = \int_{\tau} p(\tau) G(\tau) d\tau$$

其中 $p(\tau)$ 表示轨迹 τ 出现的概率。回到强化学习的目标，最大化长期奖励和的期望：

$$\arg \max_{\theta} \mathbb{E}_{\tau} J(\theta)$$

2 REINFORCE算法

2.1 目标函数的梯度

目标函数有了，一个很自然的想法就是使用梯度下降法优化策略 θ 。因此首先要做的就是求出目标函数关于 θ 的梯度。直接求导：

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \int_{\tau} p(\tau) G(\tau) d\tau \\ &= \int_{\tau} \nabla_{\theta} p(\tau) G(\tau) d\tau \\ &= \int_{\tau} p(\tau) \nabla_{\theta} \log p(\tau) G(\tau) d\tau \\ &= \mathbb{E}_{\tau} [\nabla_{\theta} \log p(\tau) G(\tau)] \end{aligned}$$

在上面的公式推导过程中使用了几个小技巧。

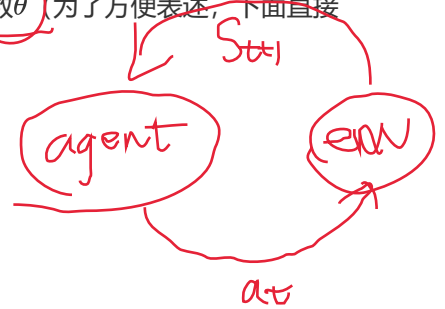
- 第1行到第2行，求导和积分交换了顺序。
- 第2行到第3行，只需要对 $p(\tau)$ 求导，而不需要对 $G(\tau)$ 求导。 θ 决定了智能体在每一步执行的动作，因而决定了整条轨迹 τ 出现的概率 $p(\tau)$ 。但在动作给出后，状态的转移以及这一步所获得奖励都是由环境的特性决定的，因此 $G(\tau)$ 不是 θ 的函数。对 $p(\tau)$ 求导的时候，使用了对数求导的特点，即 $\nabla p = p \nabla \log p$ 。

进一步考察 $\nabla_{\theta} \log p(\tau)$ ，最关键的自然是 $p(\tau)$ ，它表示整条轨迹出现的概率，和每一步的动作被选择的概率之间的关系是什么样的？因为每次做决策的时候都是独立的，因此

$$\begin{aligned} p(\tau) &= p(s_0, a_0, \dots, s_{T-1}, a_{T-1}, r_T, s_T) \\ &= p(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(s_t, a_t) p(s_{t+1} | s_t, a_t) \end{aligned}$$

进一步有

$$\begin{aligned} \nabla_{\theta} \log p(\tau) &= \nabla_{\theta} \log p(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(s_t, a_t) p(s_{t+1} | s_t, a_t) \\ &= \nabla_{\theta} \left[\log p(s_0) + \sum_{t=0}^{T-1} \log \pi_{\theta}(s_t, a_t) + \sum_{t=0}^{T-1} \log p(s_{t+1} | s_t, a_t) \right] \\ &= \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \end{aligned}$$



$$\begin{aligned} p(\tau) &= p(s_0, a_0, \dots, s_{T-1}, a_{T-1}, r_T, s_T) \\ G(\tau) &= r_1 + r_2 + \dots + r_T \end{aligned}$$

$$p(a_t | s_t) \quad p(s_{t+1} | s_t, a_t)$$

环境-智能体交互

此时目标函数的梯度可以写成

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G(\tau) \right]$$

2.2 梯度的蒙特卡洛估计

经过上面的推导后，发现梯度的表达式中存在求期望的过程。将所有可能找出来然后求期望是不可能的，只能去估计期望值，而蒙特卡洛估计则是非常有效的一种估计思想。假设我们获得了 N 条智能体与环境交互的轨迹，并且每条轨迹出现的概率相等，则目标函数梯度的蒙特卡洛估计为

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(s_t^i, a_t^i) G(\tau^i)$$

N 条轨迹 T

得到梯度后，使用梯度上升算法优化参数 θ

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

2.3 使用REINFORCE算法玩Cart-Pole

导入必要的python库，加载环境

```
1 import gym
2 gym.logger.setLevel(40) # 减少警告输出
3 import numpy as np
4 from collections import deque
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7
8 import torch
9 torch.manual_seed(500) # 随机种子
10 import torch.nn as nn
11 import torch.nn.functional as F
12 import torch.optim as optim
13 from torch.distributions import Categorical
14
15 env = gym.make('CartPole-v0')
16 env.seed(500)
17 print('观测空间:', env.observation_space)
18 print('动作空间:', env.action_space)
19
20 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
```

```
1 观测空间: Box([-4.8000002e+00 -3.4028235e+38 -4.1887903e-01 -3.4028235e+38],
2  [4.8000002e+00 3.4028235e+38 4.1887903e-01 3.4028235e+38], (4,), float32)
3 动作空间: Discrete(2)
```

定义策略神经网络

```
1 class Policy(nn.Module):
2     def __init__(self, s_size=4, h_size=16, a_size=2):
3         super(Policy, self).__init__()
4         self.fc1 = nn.Linear(s_size, h_size)
5         self.fc2 = nn.Linear(h_size, a_size)
6
7     def forward(self, x):
8         x = F.relu(self.fc1(x))
```

```

9         x = self.fc2(x)
10        return F.softmax(x, dim=1)
11
12    def act(self, state):
13        state = torch.from_numpy(state).float().unsqueeze(0).to(device)
14        probs = self.forward(state).cpu()
15        m = Categorical(probs) # 这个m里面应该包含了概率以及概率的对数
16        action = m.sample()
17        return action.item(), m.log_prob(action)

```

使用REINFORCE算法训练策略网络

```

1  policy = Policy().to(device)
2  optimizer = optim.Adam(policy.parameters(), lr=1e-2)
3
4  def reinforce(n_episodes=1000, max_t=2000, gamma=1.0, print_every=100):
5      scores_deque = deque(maxlen=100)
6      scores = []
7      for i_episode in range(1, n_episodes+1):
8          saved_log_probs = []
9          rewards = []
10         # generate a trajectory
11         state = env.reset()
12         for t in range(max_t):
13             action, log_prob = policy.act(state)
14             saved_log_probs.append(log_prob)
15             state, reward, done, _ = env.step(action)
16             rewards.append(reward)
17             if done:
18                 break
19
20         scores_deque.append(sum(rewards))
21         scores.append(sum(rewards))
22
23         discounts = [gamma**i for i in range(len(rewards))]
24         G = sum([a*b for a,b in zip(discounts, rewards)])
25
26         policy_loss = []
27         for log_prob in saved_log_probs
28             policy_loss.append(-log_prob * G) # 最大化目标，使用梯度下降，因此在目
标前加负号
29
30         policy_loss = torch.cat(policy_loss).sum() # 这个就是目标函数
31
32         optimizer.zero_grad()
33         policy_loss.backward()
34         optimizer.step()
35
36         if i_episode % print_every == 0:
37             print('Episode {} \t Average Score: {:.2f}'
38                   .format(i_episode, np.mean(scores_deque)))
39         if np.mean(scores_deque) >= 195.0:
40             print('Environment solved in {:d} episodes! \t Average Score:
41                   {:.2f}'
42                   .format(i_episode-100, np.mean(scores_deque)))
43             break

```

```

44     return scores
45
46     scores = reinforce()

```

```

1 Episode 100 Average Score: 31.32
2 Episode 200 Average Score: 52.96
3 Episode 300 Average Score: 90.69
4 Episode 400 Average Score: 168.04
5 Episode 500 Average Score: 176.96
6 Environment solved in 461 episodes! Average Score: 195.63

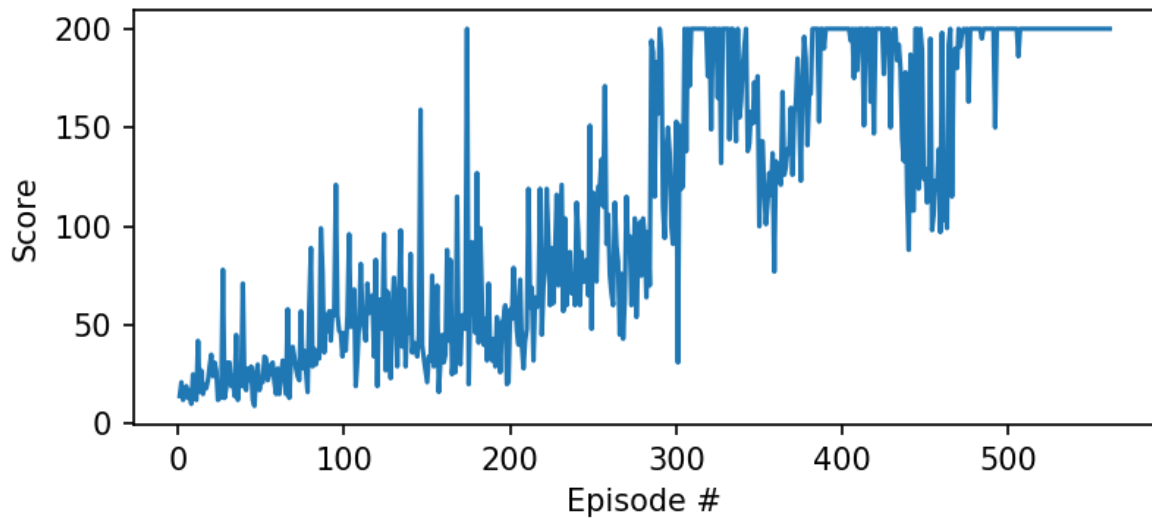
```

画训练曲线

```

1 fig = plt.figure(figsize=(6,2.5),dpi=150)
2 ax = fig.add_subplot(111)
3 plt.plot(np.arange(1, len(scores)+1), scores)
4 plt.ylabel('Score')
5 plt.xlabel('Episode #')
6 plt.show()

```



可视化训练后的智能体环境交互情况

```

1 env = gym.make('CartPole-v0')
2
3 state = env.reset()
4 for t in range(1000):
5     action, _ = policy.act(state)
6     env.render()
7     state, reward, done, _ = env.step(action)
8     if done:
9         break
10
11 env.close()

```