

# Reinforcement Learning

## Policy Gradient

袁路展

北京邮电大学

2021 年 10 月 29 日



北京邮电大学  
Beijing University of Posts and Telecommunications

① Target

② Review

③ Policy Gradient

④ Summary

# 1 Target

## 2 Review

## 3 Policy Gradient

## 4 Summary

- 什么是策略梯度？
- 为什么要策略梯度？
- 策略梯度的关键问题是？
- 一些重要的概念

REINFORCE

High variance

Score function

Actor-Critic

## ① Target

## ② Review

Review on Reinforcement Learning  
Review on Value-Based Methods

## ③ Policy Gradient

## ④ Summary

## ① Target

## ② Review

Review on Reinforcement Learning  
Review on Value-Based Methods

## ③ Policy Gradient

## ④ Summary

# Review on Reinforcement Learning

## What is Reinforcement Learning?

- Learn an optimal policy from the interact messages.

## Why Reinforcement Learning?

- High calculation complexity of classical method like DP ( $O(n^2)$ ), Optimal Control(Complex equations).
- Complex state representation.

## How to achieve the goal of RL?

- Value-based
- Policy-based
- Actor-critic

## ① Target

## ② Review

Review on Reinforcement Learning  
Review on Value-Based Methods

## ③ Policy Gradient

## ④ Summary



# Review on Value-Based Methods

- In Value-Based methods, the agent approximate the value or action-value function using parameter  $\theta$ ,

$$V_{\theta} \approx V^{\pi}(s)$$
$$Q_{\theta}(s, a) \approx Q^{\pi}(s, a)$$

- Then a policy is generated from the value function by various methods like greedy method,  $\epsilon - greedy$  and so on.
- Disadvantages of Value-based methods
  - How to select the action to take especially in continuous action space.
  - Cannot handle the stochastic policy well.

## 1 Target

## 2 Review

## 3 Policy Gradient

Intro to Policy Gradient  
Score function  
REINFORCE  
Policy Gradient Theorem

## 4 Summary

## ① Target

## ② Review

## ③ Policy Gradient

Intro to Policy Gradient

Score function

REINFORCE

Policy Gradient Theorem

## ④ Summary

# What is Policy Gradient

- 强化学习的目标是找到最优的策略  $\pi$ , 最大化收益  $J(\pi)$ .
- 基于价值函数的方法通过估计每个动作对应的价值, 然后通过选取价值最大的动作来得到最优的策略。
- 但是选取价值最大的动作在高维动作空间或者连续动作空间是非常困难的。
- 所以我们为什么不直接优化  $J(\pi)$
- 最简单的优化方法: 梯度下降

# What is Policy Gradient

- The original form of Policy Gradient:

$$\theta := \theta + \alpha \nabla J(\theta) \quad (1)$$

- Basic framework of Policy Gradient

**while**  $\theta$  *is not optimal* **do**

    | Evaluate the policy

    | Update  $\theta$  by gradient descent

**end**

- **Key problem: How to evaluate the policy and the policy gradient?**

## How to evaluate the policy

- 在强化学习中，我们的目标是让策略  $\pi_\theta$  与环境交互产生的序列  $\tau = (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$  的回报值最大。

$$J(\theta) = \mathbb{E}_{\tau \sim \rho_\theta} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t, s_{t+1}) \right] = \int_{\tau} \rho_\theta(\tau) R(\tau) d\tau \quad (2)$$

策略梯度：

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int_{\tau} \rho_{\theta}(\tau) R(\tau) d\tau = \int_{\tau} (\nabla_{\theta} \rho_{\theta}(\tau)) R(\tau) d\tau \quad (3)$$

## 1 Target

## 2 Review

## 3 Policy Gradient

Intro to Policy Gradient

**Score function**

REINFORCE

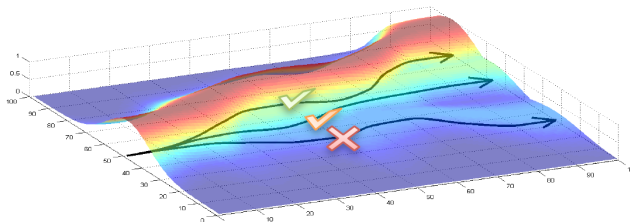
Policy Gradient Theorem

## 4 Summary

## score function

In equation(3), the gradient tries to

- Increase the probability of trajectories with higher return.
- Decrease the probability of trajectories with lower return.



**Does not try to change the trajectories!**



## score function

What's worse, because  $\rho_\theta(\tau)$  is always represented in the likelihood form:

$$\rho_\theta(\tau) = p_\theta(s_0, a_0, \dots, s_T, a_T) = p_0(s_0) \prod_{t=0}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

and the gradient of  $\prod$  is not easy to calculate.

In real analysis, log trick a common method to solve this problem.

$$\log \rho_\theta(\tau) = \log p_0(s_0) + \sum_{t=0}^T \log \pi_\theta(a_t | s_t) + \sum_{t=0}^T \log p(s_{t+1} | s_t, a_t)$$

$$\nabla_\theta \rho_\theta(\tau) = \nabla_\theta \rho_\theta(\tau) \frac{\rho_\theta(\tau)}{\rho_\theta(\tau)} = \rho_\theta(\tau) \nabla_\theta \log \rho_\theta(\tau)$$

## score function

Let's Decompose the trajectories into states and actions and the policy gradient becomes:

$$\begin{aligned}\nabla_{\theta} \log \rho_{\theta} &= \nabla_{\theta} \left[ p_0(s_0) \prod_{t=0}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t) \right] \\&= \nabla_{\theta} \left[ \log p_0(s_0) + \sum_{t=0}^T \log \pi_{\theta}(a_t | s_t) + \sum_{t=0}^T \log p(s_{t+1} | s_t, a_t) \right] \\&= \nabla_{\theta} \sum_{t=0}^T \log \pi_{\theta}(a_t | s_t) \\&= \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)\end{aligned}\tag{4}$$

**The score function is  $\nabla_{\theta} \log \pi_{\theta}(a|s)$**

## score function

The policy gradient becomes:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int_{\tau} \rho_{\theta}(\tau) \nabla_{\theta} \log \rho_{\theta}(\tau) R(\tau) d\tau \\&= \mathbb{E}_{\tau \sim \rho_{\theta}} [\nabla_{\theta} \log \rho_{\theta} R(\tau)] \\&= \mathbb{E}_{\tau \sim \rho_{\theta}} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) R(\tau) \right] \\&\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \quad (5)\end{aligned}$$

## advantage of score function

- No dynamics in the objective function.
- Focus on the policy instead of the trajectory.
- Better computing feathers.
  - Most stochastic policies are represented in exponential form.

## 1 Target

## 2 Review

## 3 Policy Gradient

Intro to Policy Gradient

Score function

**REINFORCE**

Policy Gradient Theorem

## 4 Summary

# Monte-Carlo Policy(REINFORCE)

- REINFORCE

Initialise  $\theta$

**for** *each episode* **do**

    Generate an episode  $\tau^i$  by  $\pi(\theta)$

$$\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$$

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

**end**

return  $\theta$

- **pay attention to**  $v_t$

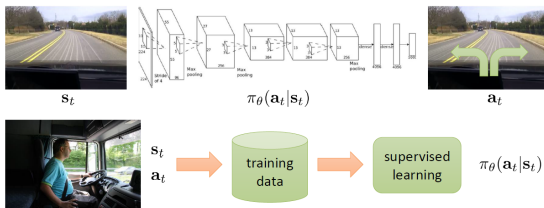
# Compaision to maximum likelihood

policy gradient

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \quad (6)$$

maximum likelihood

$$\nabla_{\theta} J_{\text{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \quad (7)$$



# Restriction on REINFORCE

While very simple, REINFORCE does not work well in practice:

- The return  $R(\tau)$  have a very high variance.

- Sensitive to the reward.

- It requires a lot of episodes to converge.

- It only works with online learning.

- Restricted in episodic environments.



# High Variance in REINFORCE

- 想象一下，你在玩英雄联盟，你刚开始进入游戏采取的动作都是一样的，但是有的局你赢了，有的局你输了，那么你最开始的动作，应该向哪个方向优化？
- 考虑两个环境 A 和 B，他们有相同的 dynamics 和相同的任务，但是 reward 设置不一样，在 A 中，执行动作 1 的奖励是 1000，执行动作 2 的奖励是 1001，在 B 中执行动作 1 的奖励是 0，执行动作 2 的奖励是 1，这两种情况下，显然 A 环境下收敛速度要缓慢。
- 在有监督学习中，high variance 的问题很少，因为训练集是固定的，但是在强化学习中，你不可能进行足够多的采样来覆盖整个动作-序列空间。

# High Variance in REINFORCE

The problem is even worse in the following conditions:

- High-dimensional action spaces: it becomes difficult to sample the environment densely enough if many actions are possible.
- Long horizons: the longer the trajectory, the more likely it will be unique.
- Finite samples: if we cannot sample enough trajectories, the high variance can introduce a bias in the gradient, leading to poor convergence.

## 1 Target

## 2 Review

## 3 Policy Gradient

Intro to Policy Gradient  
Score function  
REINFORCE  
Policy Gradient Theorem

## 4 Summary

# Policy Gradient Theorem

## 定理

### *Policy Gradient Theorem*

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{S}} \rho^{\pi_{\theta}}(s) \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) Q^{\pi}(s, a) da ds$$

- Because the actual return  $R(\tau)$  is replaced by its expectation  $Q(s,a)$ , the policy gradient is now mathematical expectation over single transition instead of the complete trajectories.
- $\rho^{\pi_{\theta}}(s)$  depends on  $\theta$ , but there is no  $\nabla_{\theta} \rho^{\pi_{\theta}}(s)$  term in  $\nabla_{\theta} J(\theta)$

- ① Target
- ② Review
- ③ Policy Gradient
- ④ Summary

# Why Policy Gradient

## Advantages of PG

- Learn the final object of RL directly. (*learn the optimal policy directly*)
- Easy extend to **high-dimensional or continous** action sapce. (*no argmax operator*)
- Can learn **stochastic** policies.

## Disadvantages of PG

- Get stuck in local optimal. (*common problem of gradient methods*)
- Evaluating a policy is typically inefficient and **high variance**.

## Review on this class

- 什么是策略梯度？
  - 直接根据  $J(\pi)$  的梯度更新策略的参数。

为什么要策略梯度？

- 基于值函数的方法不能很好应对高维、连续动作空间。

策略梯度的关键问题是什么？

- 策略的评估!!!!

# Review on this class

关键知识点:

- Score function
  - 为什么用 log? 对 trajectory 的分解。

REINFORCE

- Monte-Carlo policy gradient.

Policy Gradient Theorem

- 不再是对整个路径的积分, 只考虑单个 **transition!!** 为后面 **Actor-Critic** 结构奠定基础



# Next Class

- Advantage Policy Gradient (*reduce the variance of policy gradient*)
- Actor-Critic (*Most common algorithm frame in RL*)
- off-policy Actor-Critic (*practical policy evaluation*)

# Explore yourself

- TRPO, PPO. (*montonic improvement based in policy gradient*)
- DPG, DDPG, TD3. (*introduce success tricks in DQN to PG*)
- soft-RL, SAC, entropy-based methods. (*Explore-Exploit in PG*)
- off-policy off-line policy evaluation.

# Questions

## Q&A

