

强化学习2024

第8节

涉及知识点：

基于神经网络的策略梯度、A3C、确定性梯度策略、深度确定性策略梯度、TRPO、PPO



深度策略梯度

张伟楠 - [上海交通大学](#)

2024年上海交通大学ACM班强化学习课程大纲

强化学习基础部分

(中文课件)

1. 强化学习、探索与利用
2. MDP和动态规划
3. 值函数估计
4. 无模型控制方法
5. 规划与学习
6. 参数化的值函数和策略
7. 深度强化学习价值方法
8. 深度强化学习策略方法

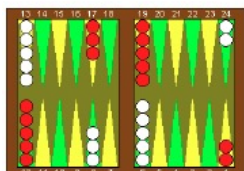
强化学习前沿部分

(英文课件)

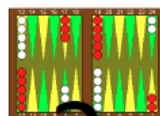
9. 基于模型的深度强化学习
10. 离线强化学习
11. 模仿学习
12. 多智能体强化学习基础
13. 多智能体强化学习前沿
14. 基于扩散模型的强化学习
15. AI Agent与决策大模型
16. 技术与交流与回顾

深度强化学习

标准 (传统)
强化学习



features



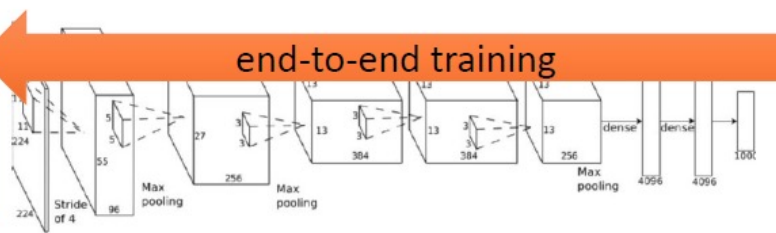
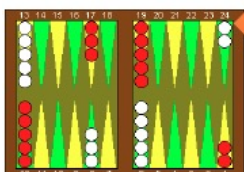
more features



linear policy
or value func.

action

深度强化学习



action

- 直面理解：深度学习+强化学习
- 深度强化学习使强化学习算法能够以端到端的方式解决复杂问题
- 真正让强化学习有能力完成实际决策任务
- 比强化学习和深度学习各自都更加难以驯化
- 基于价值函数的深度强化学习
 - DQN：一次输入多个行动Q值输出、目标网络、随机采样经验
 - Double DQN：解耦合行动选择和价值估计、解决DQN过高估计问题
 - Dueling DQN：精细捕捉价值和行动的细微关联、多种advantage函数建模

深度强化学习的分类

□ 基于价值的方法

- 深度Q网络及其扩展

□ 基于确定性策略的方法

- 确定性策略梯度 (DPG) , DDPG , TD3

□ 基于随机策略的方法

- 使用神经网络的策略梯度 , 自然策略梯度 , 信任区域策略优化 (TRPO) , 近端策略优化 (PPO) , A3C



基于神经网络的策略梯度

讲师：张伟楠 - [上海交通大学](#)

复习：策略梯度定理

- 策略梯度定理把似然比的推导过程泛化到多步马尔可夫决策过程
 - 用长期的价值函数 $Q^{\pi_\theta}(s, a)$ 代替前面的瞬时奖励 r_{sa}
- 策略梯度定理涉及
 - 起始状态目标函数 J_1 , 平均奖励目标函数 J_{avR} , 和平均价值目标函数 J_{avV}
- 定理
 - 对任意可微的策略 $\pi_\theta(a|s)$, 任意策略的目标函数 $J = J_1, J_{avR}, J_{avV}$, 其策略梯度是

$$\frac{\partial J(\theta)}{\partial \theta} = \mathbb{E}_{\pi_\theta} \left[\frac{\partial \log \pi_\theta(a|s)}{\partial \theta} Q^{\pi_\theta}(s, a) \right]$$

详细证明过程请参考 Rich Sutton's Reinforcement Learning: An Introduction (2nd Edition) 第13章

策略网络的梯度

- 对于随机策略，一般采样到每一个行动的概率由Softmax实现

$$\pi_{\theta}(a|s) = \frac{e^{f_{\theta}(s,a)}}{\sum_{a'} e^{f_{\theta}(s,a')}}$$

- 其中 $f_{\theta}(s, a)$ 是对状态-行动队的打分函数，由 θ 参数化，这可以通过一个神经网络来实现

- 其log形式的梯度为

$$\begin{aligned} \frac{\partial \log \pi_{\theta}(a|s)}{\partial \theta} &= \frac{\partial f_{\theta}(s, a)}{\partial \theta} - \frac{1}{\sum_{a'} e^{f_{\theta}(s,a')}} \sum_{a''} e^{f_{\theta}(s,a'')} \frac{\partial f_{\theta}(s, a'')}{\partial \theta} \\ &= \frac{\partial f_{\theta}(s, a)}{\partial \theta} - \mathbb{E}_{a' \sim \pi_{\theta}(a'|s)} \left[\frac{\partial f_{\theta}(s, a')}{\partial \theta} \right] \end{aligned}$$

策略网络的梯度

□ 其log梯度形式

$$\frac{\partial \log \pi_{\theta}(a|s)}{\partial \theta} = \frac{\partial f_{\theta}(s, a)}{\partial \theta} - \mathbb{E}_{a' \sim \pi_{\theta}(a'|s)} \left[\frac{\partial f_{\theta}(s, a')}{\partial \theta} \right]$$

□ 策略网络的梯度为

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \mathbb{E}_{\pi_{\theta}} \left[\frac{\partial \log \pi_{\theta}(a|s)}{\partial \theta} Q^{\pi_{\theta}}(s, a) \right] \\ &= \mathbb{E}_{\pi_{\theta}} \left[\underbrace{\left(\frac{\partial f_{\theta}(s, a)}{\partial \theta} \right)}_{\text{反向梯度传播}} - \mathbb{E}_{a' \sim \pi_{\theta}(a'|s)} \left[\underbrace{\left(\frac{\partial f_{\theta}(s, a')}{\partial \theta} \right)}_{\text{反向梯度传播}} \right] Q^{\pi_{\theta}}(s, a) \right] \end{aligned}$$

策略梯度和Q学习的对比

- Q 学习算法学习一个由 θ 作为参数的函数 $Q_\theta(s, a)$

- 优化目标为最小化TD error

$$J(\theta) = \mathbb{E}_{\pi'} \left[\frac{1}{2} \left(r_t + \gamma \max_{a'} Q_{\theta'}(s_{t+1}, a') - Q_\theta(s_t, a_t) \right)^2 \right]$$

- 更新方程

$$\begin{aligned} \theta &\leftarrow \theta - \alpha \frac{\partial J(\theta)}{\partial \theta} \\ &= \theta + \alpha \mathbb{E}_{\pi'} \left[\left(r_t + \gamma \max_{a'} Q_{\theta'}(s_{t+1}, a') - Q_\theta(s_t, a_t) \right) \frac{\partial Q_\theta(s, a)}{\partial \theta} \right] \end{aligned}$$

- 策略梯度学习一个由 θ 作为参数的策略 $\pi_\theta(a|s)$

- 优化目标直接为策略的价值 (比Q学习更加直接)

$$\max_{\theta} J(\theta) = \mathbb{E}_{\pi_\theta} [A^{\pi_\theta}(s, a)]$$

- 更新方程

$$\theta \leftarrow \theta + \alpha \frac{\partial J(\theta)}{\partial \theta} = \theta + \alpha \mathbb{E}_{\pi_\theta} \left[\frac{\partial \log \pi_\theta(a|s)}{\partial \theta} A^{\pi_\theta}(s, a) \right]$$



A3C

讲师：张伟楠 - [上海交通大学](#)

复习：Actor-Critic

Actor-Critic的思想

- REINFORCE策略梯度方法：使用蒙特卡罗采样直接估计 (s_t, a_t) 的值 G_t
- 为什么不建立一个可训练的值函数 Q_ϕ 来完成这个估计过程？

演员 (Actor) 和评论家 (Critic)

演员 $\pi_\theta(a|s)$

采取动作使评论家满意的策略



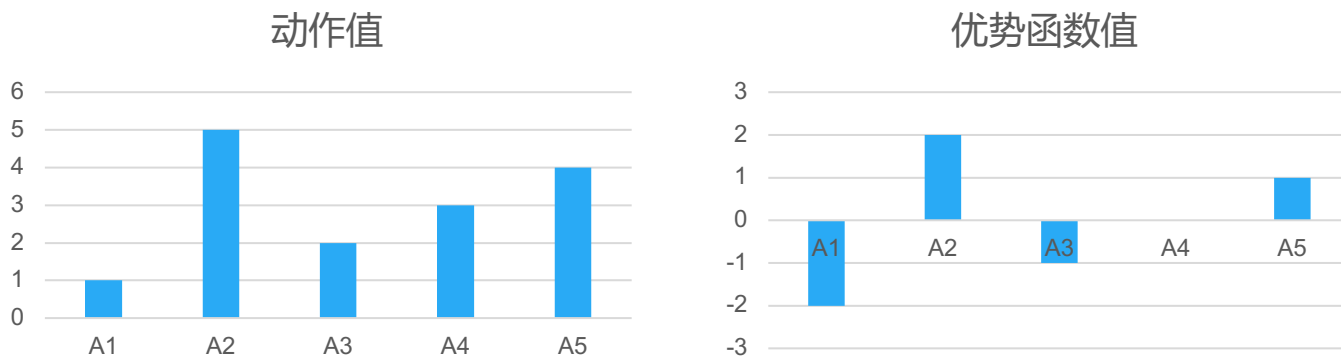
评论家 $Q_\phi(s, a)$

学会准确估计演员策略所采取动作价值的值函数

复习A2C：优势Actor-Critic

- 思想：通过减去一个基线函数来标准化评论家的打分
 - 更多信息指导：降低较差动作概率，提高较优动作概率
 - 进一步降低方差
- 优势函数 (Advantage Function)

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$



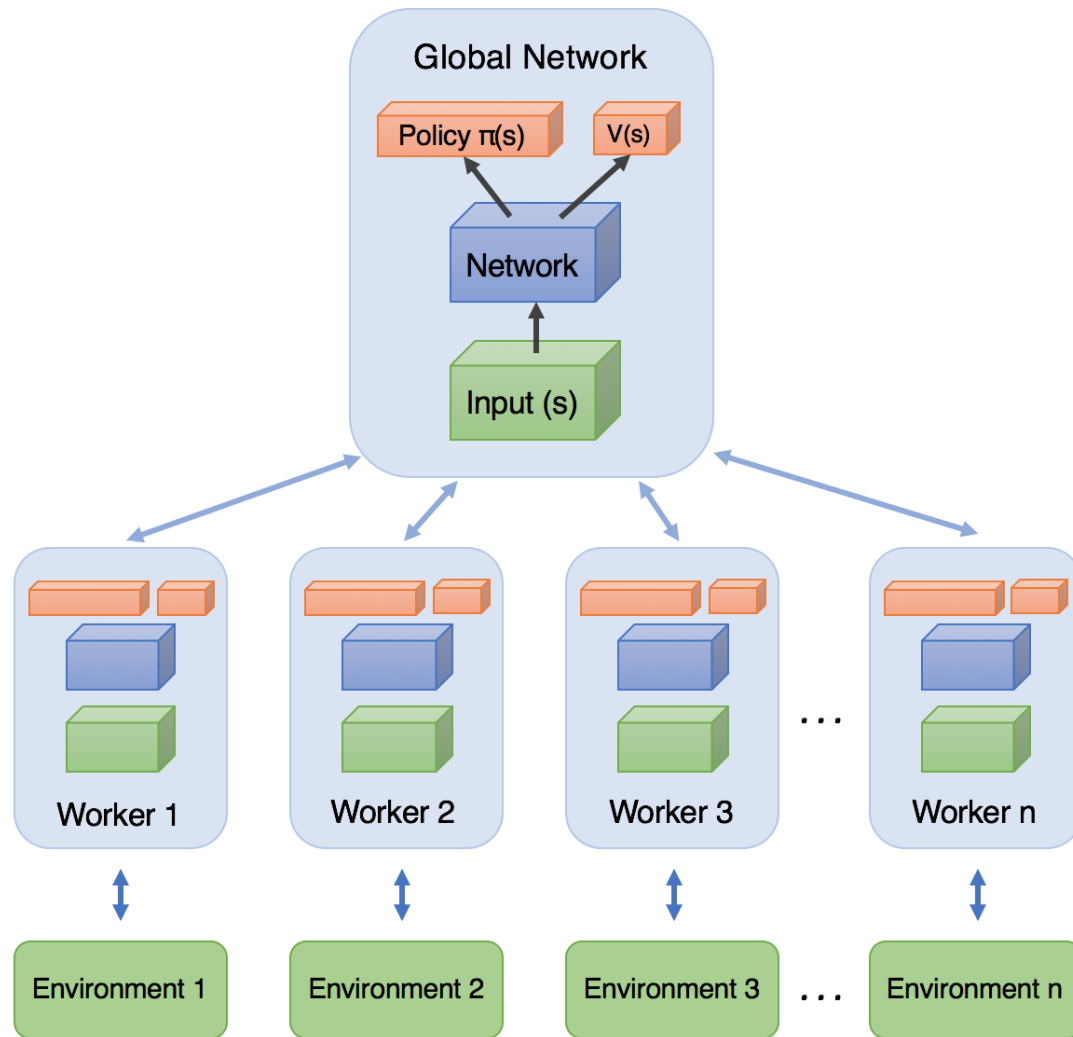
A3C: 异步A2C方法

- A3C代表了异步优势动作评价 (Asynchronous Advantage Actor Critic)
 - 异步 (Asynchronous) : 因为算法涉及并行执行一组环境
 - 优势 (Advantage) : 因为策略梯度的更新使用优势函数
 - 动作评价 (Actor Critic) : 因为这是一种动作评价 (actor-critic) 方法 , 它涉及一个在学得的状态值函数帮助下进行更新的策略

$$\nabla_{\theta'} \log \pi(a_t | s_t; \theta') A(s_t, a_t; \theta_v)$$

$$A(s_t, a_t; \theta_v) = \sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v)$$

A3C: 异步A2C方法



A3C算法

Algorithm S3 Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

// Assume global shared parameter vectors θ and θ_v and global shared counter $T = 0$

// Assume thread-specific parameter vectors θ' and θ'_v

Initialize thread step counter $t \leftarrow 1$

repeat

Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.

Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$

$t_{start} = t$

Get state s_t

repeat

Perform a_t according to policy $\pi(a_t|s_t; \theta')$

Receive reward r_t and new state s_{t+1}

$t \leftarrow t + 1$

$T \leftarrow T + 1$

until terminal s_t **or** $t - t_{start} == t_{max}$

$R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{ Bootstrap from last state} \end{cases}$

for $i \in \{t - 1, \dots, t_{start}\}$ **do**

$R \leftarrow r_i + \gamma R$

Accumulate gradients wrt θ' : $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$

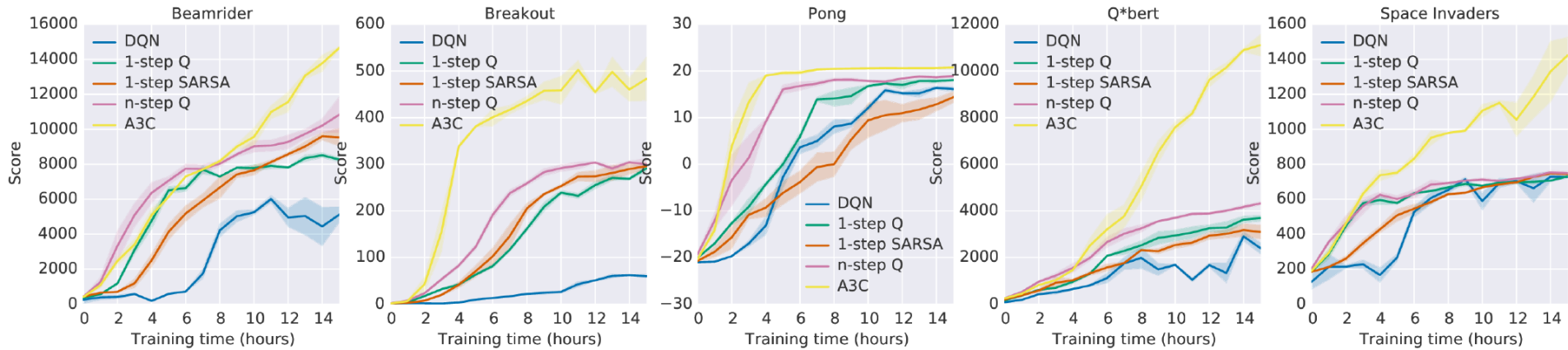
Accumulate gradients wrt θ'_v : $d\theta_v \leftarrow d\theta_v + \partial (R - V(s_i; \theta'_v))^2 / \partial \theta'_v$

end for

Perform asynchronous update of θ using $d\theta$ and of θ_v using $d\theta_v$.

until $T > T_{max}$

A3C对比实验



a single Nvidia K40 GPU while the asynchronous methods were trained using 16 CPU cores

Method	Training Time	Mean	Median
DQN	8 days on GPU	121.9%	47.5%
Gorila	4 days, 100 machines	215.2%	71.3%
D-DQN	8 days on GPU	332.9%	110.9%
Dueling D-DQN	8 days on GPU	343.8%	117.1%
Prioritized DQN	8 days on GPU	463.6%	127.6%
A3C, FF	1 day on CPU	344.1%	68.2%
A3C, FF	4 days on CPU	496.8%	116.6%
A3C, LSTM	4 days on CPU	623.0%	112.6%

} Nvidia K40 GPUs
 } 16 CPU cores and no GPU

Mean and median human-normalized scores on 57 Atari games



信任区域策略优化 TRPO

讲师：张伟楠 - [上海交通大学](#)

目录

Contents

- 01 策略梯度的缺点
- 02 TRPO算法
- 03 策略改进的单调性保证
- 04 实验结果

01

策略梯度的缺点

策略梯度算法回顾

□ 蒙特卡洛策略梯度 (REINFORCE) 算法

initialize θ arbitrarily

for each episode $\{s_1, a_1, r(s_1, a_1), \dots, s_T, a_T, r(s_T, a_T)\} \sim \pi_\theta$ do

for $t = 1$ to T do

$$\theta \leftarrow \theta + \alpha \frac{\partial}{\partial \theta} \log \pi_\theta(a_t | s_t) G_t$$

end for

end for

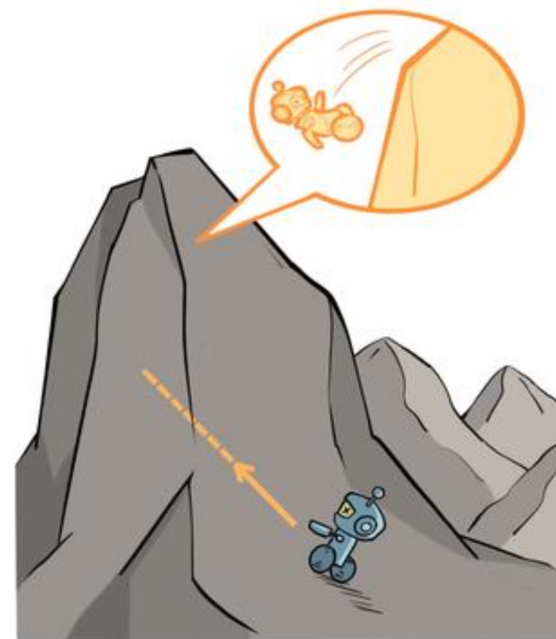
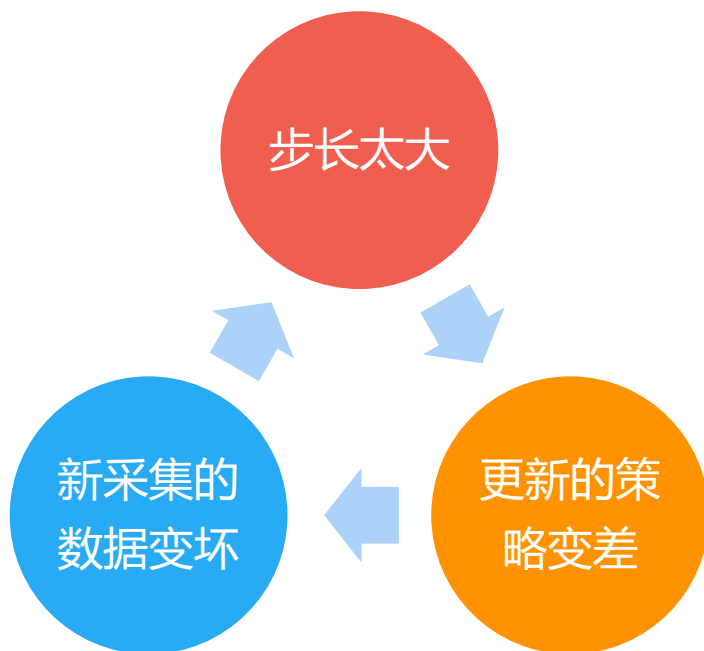
return θ

相关定义

- $s_t, a_t, r(s_t, a_t)$: t 时刻的状态, 动作和奖励
- π_θ, θ : 使用的策略, 表示策略所使用的参数
- G_t : 累计奖励
- α : 步长

策略梯度的缺点

- 适合的步长在策略梯度中难以确定
 - 采集到的数据的分布会随策略的更新而变化。
 - 较差的步长产生的影响大。





策略梯度的优化目标

□ 优化目标的两种形式

- 第一种： $J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\sum_t \gamma^t r(s_t, a_t)]$
- 因为 $V^{\pi_{\theta}}(s) = \mathbb{E}_{a \sim \pi_{\theta}(s)} [Q^{\pi_{\theta}}(s, a)] = \mathbb{E}_{a \sim \pi_{\theta}(s)} [\mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\sum_{s_k=s, a_k=a} \sum_{t=k}^{\infty} \gamma^{t-k} r(s_t, a_t)]]$ 。
- 所以优化目标的第二种形式是： $J(\theta) = \mathbb{E}_{s_0 \sim p_{\theta}(s_0)} [V^{\pi_{\theta}}(s_0)]$

相关定义

- τ ：轨迹
- s_0 ：初始状态
- $s_t, a_t, r(s_t, a_t)$ ： t 时刻的状态，动作和奖励
- π_{θ} ：使用的策略
- θ ：表示策略所使用的参数
- $Q^{\pi_{\theta}}$ 和 $V^{\pi_{\theta}}$ ：策略 π_{θ} 下的 Q 值与状态值函数

优化目标的优化量

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\sum_t \gamma^t r(s_t, a_t)]$$
$$J(\theta) = \mathbb{E}_{s_0 \sim p_{\theta}(s_0)} [V^{\pi_{\theta}}(s_0)]$$

初始状态的分布与 θ 无关

$$\begin{aligned} J(\theta') - J(\theta) &= J(\theta') - \mathbb{E}_{s_0 \sim p(s_0)} [V^{\pi_{\theta}}(s_0)] \\ &= J(\theta') - \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} [V^{\pi_{\theta}}(s_0)] \\ &= J(\theta') - \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t V^{\pi_{\theta}}(s_t) - \sum_{t=1}^{\infty} \gamma^t V^{\pi_{\theta}}(s_t) \right] \\ &= J(\theta') + \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_{\theta}}(s_{t+1}) - V^{\pi_{\theta}}(s_t)) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] + \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_{\theta}}(s_{t+1}) - V^{\pi_{\theta}}(s_t)) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \gamma V^{\pi_{\theta}}(s_{t+1}) - V^{\pi_{\theta}}(s_t)) \right] \\ &= \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right] \end{aligned}$$

$J(\theta')$ 的定义

不方便采样

$$A^{\pi_{\theta}}(s_t, a_t) = Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)$$

使用重要性采样

□ 使用重要性采样 (Importance Sampling)

$$\begin{aligned} & J(\theta') - J(\theta) \\ &= \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta'}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta'}(a_t|s_t)} [\gamma^t A^{\pi_{\theta}}(s_t, a_t)]] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta'}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} \left[\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]] \end{aligned}$$

$$\begin{aligned} & A^{\pi_{\theta}}(s_t, a_t) \\ &= Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t) \end{aligned}$$

仍然是 $p_{\theta'}$
(近似操作)

重要性采样

忽略状态分布的差异

- 当策略更新前后的变化较小时，可以令 $p_\theta(s_t) \approx p_{\theta'}(s_t)$ 。
 - 假设使用确定性策略，当 $\pi_{\theta'}(s_t) \neq \pi_\theta(s_t)$ 的概率小于 ϵ 时
 - 或者假设使用随机策略，当 $a' \sim \pi_{\theta'}(\cdot | s_t) \neq a \sim \pi_\theta(\cdot | s_t)$ 的概率小于 ϵ 时
 - $p_{\theta'}(s_t) = (1 - \epsilon)^t p_\theta(s_t) + (1 - (1 - \epsilon)^t) p_{\text{mistake}}(s_t)$
 - $|p_{\theta'}(s_t) - p_\theta(s_t)| = (1 - (1 - \epsilon)^t) |p_{\text{mistake}}(s_t) - p_\theta(s_t)| \leq 2(1 - (1 - \epsilon)^t) \leq 2\epsilon t$

$$(1 - \epsilon)^t \geq 1 - \epsilon t \text{ for } \epsilon \in [0,1]$$

$$J(\theta') - J(\theta) \approx \sum_t \mathbb{E}_{s_t \sim p_\theta(s_t)} [\mathbb{E}_{a_t \sim \pi_\theta(a_t | s_t)} [\frac{\pi_{\theta'}(a_t | s_t)}{\pi_\theta(a_t | s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t)]]$$

约束策略的变化

- 使用KL散度约束策略更新的幅度

$$\theta' \leftarrow \arg \max_{\theta'} \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} [\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]]$$

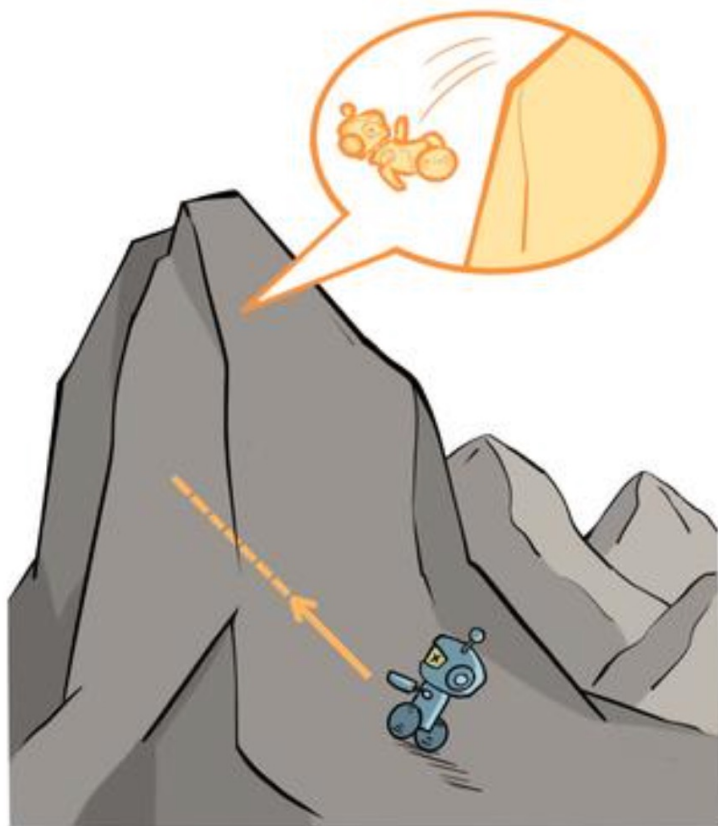
such that $\mathbb{E}_{s_t \sim p_{\theta}(s_t)} [D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t))] \leq \epsilon$

- 可以使用基于泰勒展开和共轭梯度的求解方法（见动手学强化学习11.3节）
- 实际多使用constraint violate as penalty

$$\theta' \leftarrow \arg \max_{\theta'} \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} [\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]] - \lambda(D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) - \epsilon)$$

1. 优化上式，更新 θ'
2. 更新 $\lambda \leftarrow \lambda + \alpha(D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) - \epsilon)$

TRPO的原理



Gradient Ascent



Optimization in Trust Region

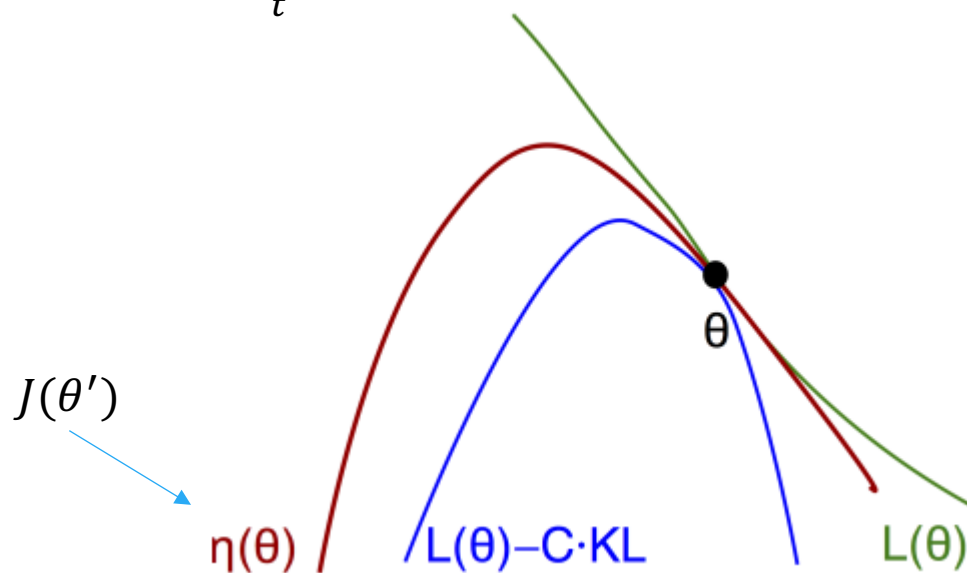
03

策略改进的
单调性保证

策略改进的单调性保证

$$J(\theta') \geq L_{\theta}(\theta') - C \cdot D_{KL}^{max}(\theta, \theta'), \text{ where } C = \frac{4\epsilon\gamma}{(1-\gamma)^2}, \epsilon = \max_{s,a} |A_{\pi}(s, a)|$$

$$L_{\theta}(\theta') = J(\theta) + \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} \left[\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]]$$

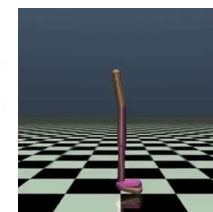
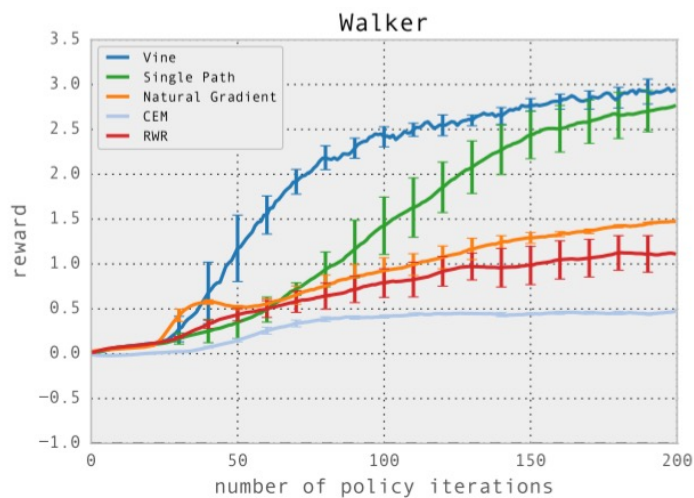
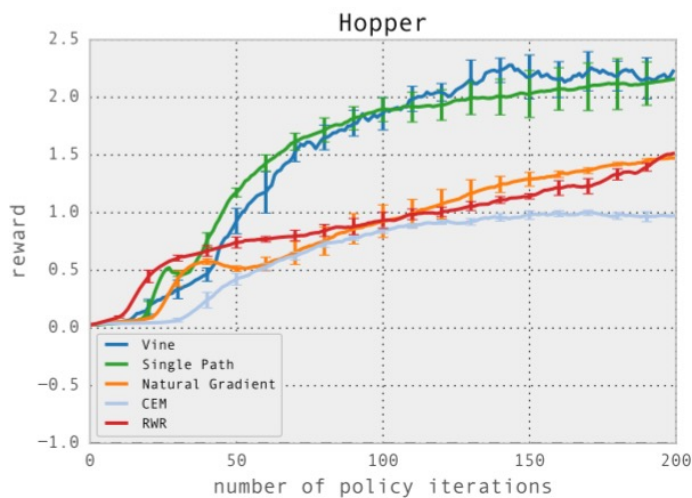
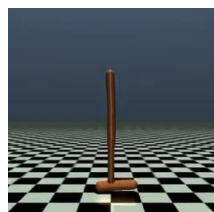
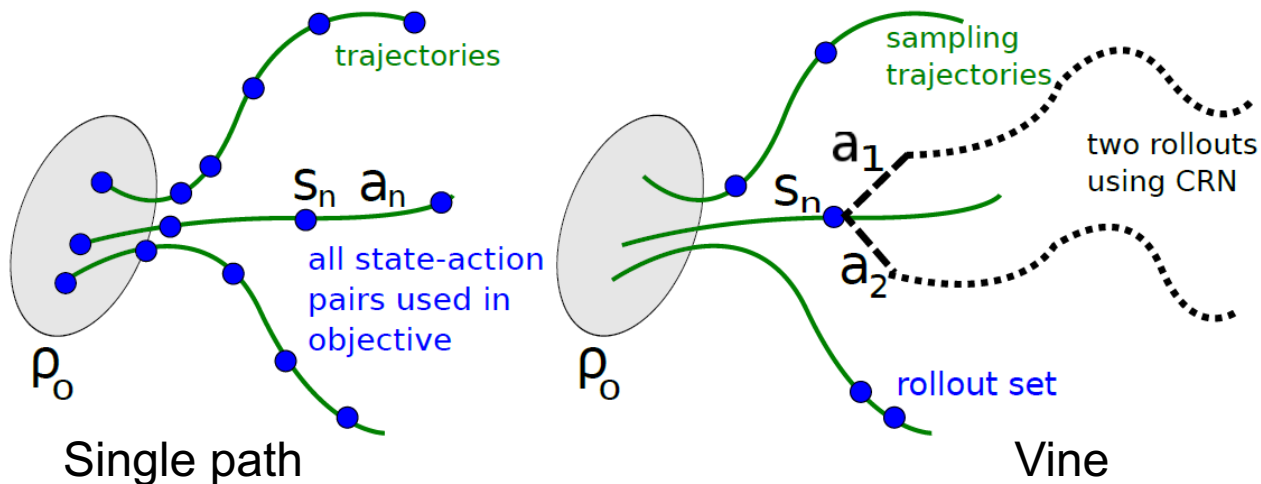




04

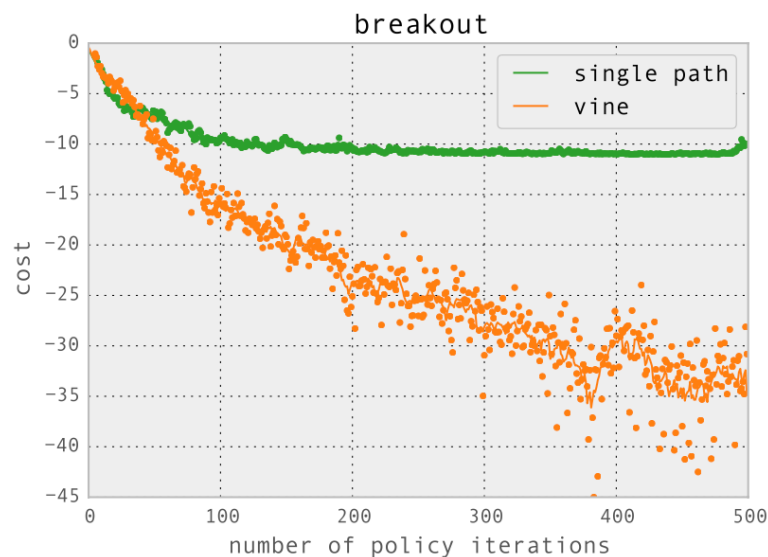
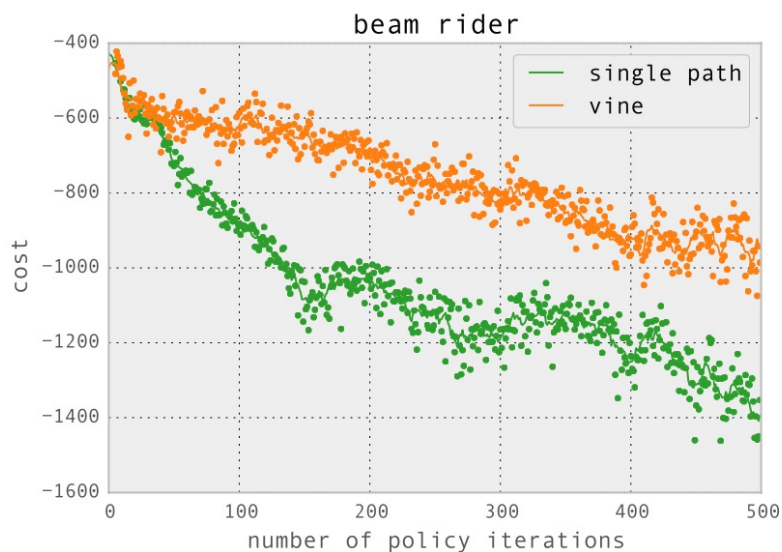
实验结果

训练曲线



结果比较

	<i>B. Rider</i>	<i>Breakout</i>	<i>Enduro</i>	<i>Pong</i>	<i>Q*bert</i>	<i>Seaquest</i>	<i>S. Invaders</i>
Random	354	1.2	0	-20.4	157	110	179
Human (Mnih et al., 2013)	7456	31.0	368	-3.0	18900	28010	3690
Deep Q Learning (Mnih et al., 2013)	4092	168.0	470	20.0	1952	1705	581
UCC-I (Guo et al., 2014)	5702	380	741	21	20025	2995	692
TRPO - single path	1425.2	10.8	534.6	20.9	1973.5	1908.6	568.4
TRPO - vine	859.5	34.2	430.8	20.9	7732.5	788.4	450.2



❑ TRPO的不足

- 近似带来误差
- 求解约束优化问题的困难

❑ PPO算法

- 理论更简洁，操作更简单，实验效果更好
- 推荐阅读 Proximal Policy Optimization Algorithms, [John Schulman](#), et al. (2017)



John Schulman

Research Scientist, OpenAI
 Verified email at openai.com - [Homepage](#)
[Artificial Intelligence](#) [Robotics](#) [Neuroscience](#)

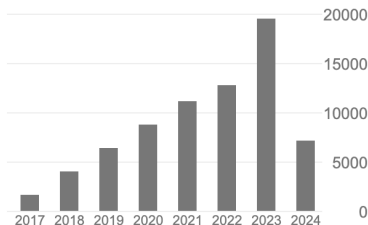
[FOLLOW](#)

TITLE	CITED BY	YEAR
Proximal policy optimization algorithms J Schulman, F Wolski, P Dhariwal, A Radford, O Klimov arXiv preprint arXiv:1707.06347	16893	2017
Trust region policy optimization J Schulman, S Levine, P Abbeel, M Jordan, P Moritz International conference on machine learning, 1889-1897	7815	2015
OpenAI Gym G Brockman, V Cheung, L Pettersson, J Schneider, J Schulman, J Tang, ... arXiv preprint arXiv:1606.01540	7121	2016
Training language models to follow instructions with human feedback L Ouyang, J Wu, X Jiang, D Almeida, C Wainwright, P Mishkin, C Zhang, ... Advances in neural information processing systems 35, 27730-27744	5970	2022

Google Scholar

Cited by

	All	Since 2019
Citations	72863	66022
h-index	54	53
i10-index	72	70



Public access

[VIEW ALL](#)

0 articles

7 articles



近端策略优化

Proximal Policy Optimization

讲师：张伟楠 - [上海交通大学](#)

回顾TRPO

- TRPO使用KL散度约束策略更新的幅度

$$\theta' \leftarrow \arg \max_{\theta'} \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} [\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]]$$

such that $D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) \leq \epsilon$

- 使用constraint violate as penalty

$$\theta' \leftarrow \arg \max_{\theta'} \sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [\mathbb{E}_{a_t \sim \pi_{\theta}(a_t|s_t)} [\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t)]] - \lambda(D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) - \epsilon)$$

1. 优化上式，更新 θ'
2. 更新 $\lambda \leftarrow \lambda + \alpha(D_{KL}(\pi_{\theta'}(a_t|s_t) \parallel \pi_{\theta}(a_t|s_t)) - \epsilon)$

TRPO的不足

- 重要性比例带来的大方差
- 求解约束优化问题的困难

PPO: Proximal Policy Optimization

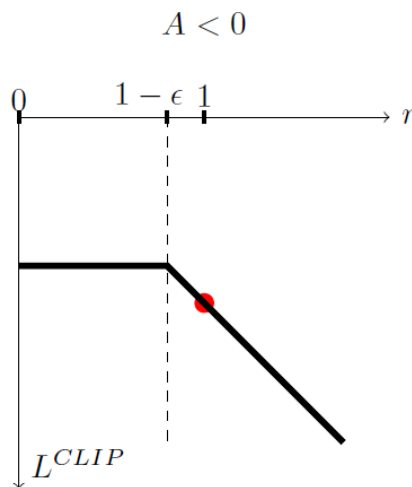
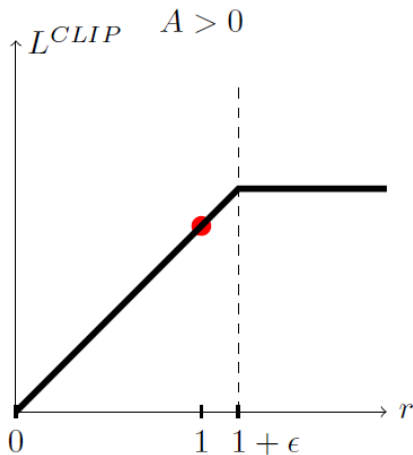
PPO在TRPO基础上的改进

1. 截断式优化目标

conservative
policy iteration

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)]$$



构建下界

$$L^{CLIP}(\theta) \leq L^{CPI}(\theta)$$

在 $r = 1$ 附近相等

$$L^{CLIP}(\theta) = L^{CPI}(\theta)$$

PPO: Proximal Policy Optimization

PPO在TRPO基础上的改进

1. 截断式优化目标

conservative
policy iteration

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)]$$

2. 优势函数 \hat{A}_t 选用多步时序差分

$$\hat{A}_t = -V(s_t) + r_t + \gamma r_{t+1} + \dots + \gamma^{T-t+1} r_{T-1} + \gamma^{T-t} V(s_T)$$

- 在每次迭代中，并行 N 个actor收集 T 步经验数据
- 计算每步的 \hat{A}_t 和 $L^{CLIP}(\theta)$ ，构成mini-batch
- 更新参数 θ ，并更新 $\theta_{\text{old}} \leftarrow \theta$

PPO: Proximal Policy Optimization

PPO在TRPO基础上的改进

3. 自适应的KL惩罚项参数

$$L^{KL PEN}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t - \beta \text{KL}[\pi_{\theta_{old}}(\cdot | s_t) | \pi_{\theta}(\cdot | s_t)] \right]$$

动态调整 β 方法

- 计算KL值 $d = \hat{\mathbb{E}}_t \left[\text{KL}[\pi_{\theta_{old}}(\cdot | s_t) | \pi_{\theta}(\cdot | s_t)] \right]$
 - a) 如果 $d < d_{\text{targ}}/1.5$, 更新 $\beta \leftarrow \beta/2$
 - b) 如果 $d > d_{\text{targ}} \times 1.5$, 更新 $\beta \leftarrow \beta \times 2$

注：这里1.5和2是经验参数，算法效能和它们并不是很敏感

PPO实验对比

No clipping or penalty:

$$L_t(\theta) = r_t(\theta)\hat{A}_t$$

Clipping:

$$L_t(\theta) = \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta)), 1 - \epsilon, 1 + \epsilon)\hat{A}_t$$

KL penalty (fixed or adaptive)

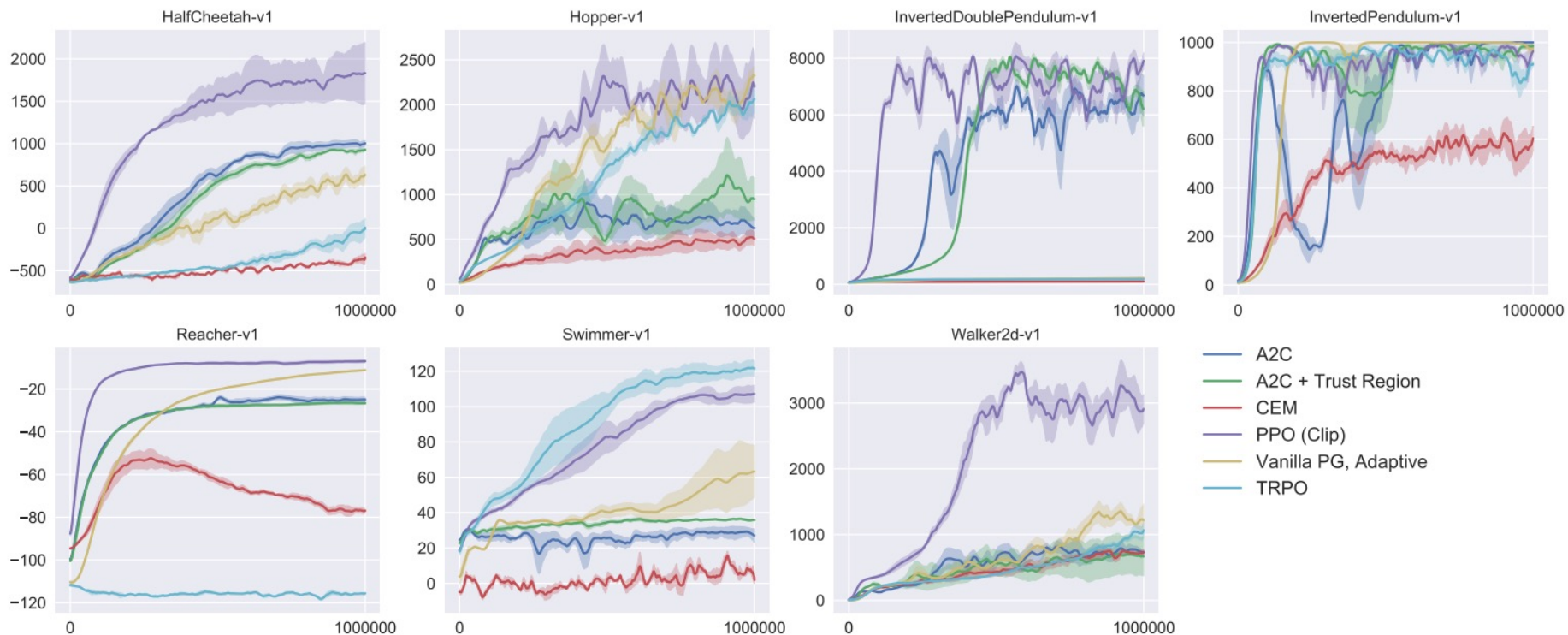
$$L_t(\theta) = r_t(\theta)\hat{A}_t - \beta \text{KL}[\pi_{\theta_{\text{old}}}, \pi_{\theta}]$$

- 7个连续控制的环境
- 3个random seed
- 每个算法跑100个 episode , 跑21遍 , 做平均值计算
- 最佳score归一化为1

algorithm	avg. normalized score
No clipping or penalty	-0.39
Clipping, $\epsilon = 0.1$	0.76
Clipping, $\epsilon = 0.2$	0.82
Clipping, $\epsilon = 0.3$	0.70
Adaptive KL $d_{\text{targ}} = 0.003$	0.68
Adaptive KL $d_{\text{targ}} = 0.01$	0.74
Adaptive KL $d_{\text{targ}} = 0.03$	0.71
Fixed KL, $\beta = 0.3$	0.62
Fixed KL, $\beta = 1.$	0.71
Fixed KL, $\beta = 3.$	0.72
Fixed KL, $\beta = 10.$	0.69

PPO实验对比

- 在MuJoCo实验环境中的对比结果（训练100万步）



PPO在ChatGPT中的使用

加入了基于人类的反馈系统

Reinforcement Learning from Human Feedback

丛问题库里抽取问题

什么是香蕉?

标记者 (Labeler) 书写期待的回复

香蕉是一种水果, 从香蕉树...

被标记的数据用来调优 GPT-3.5



采样问题, 并列出所有模型和标记者的回答

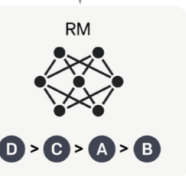
什么是香蕉?

- A 香蕉是一种水果, 从香蕉树...
- B 香蕉是芭蕉科、芭蕉属植物...
- C 香蕉, 从属性来说, 与草莓、葡萄、猕猴桃同类...
- D 香蕉为芭蕉科植物甘蕉的果实, 原产亚洲东南部...

标记者 (Labeler) 排序所有标记着答案

D > C > A > B

用排序答案训练奖励模型



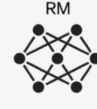
通过模型生成初步回答

写个水獭的故事



很久很久以前...

输入奖励模型得到分数和优化参数



r_k

持续优化参数迭代

总结深度策略梯度方法

- 相比价值函数学习最小化TD误差的目标，策略梯度方法直接优化策略价值的目标更加贴合强化学习本质目标
- 分布式的actor-critic算法能够充分利用多核CPU资源采样环境的经验数据，利用GPU资源异步地更新网络，这有效提升了DRL的训练效率
- 基于神经网络的策略在优化时容易因为一步走得太大而变得很差，进而下一轮产生很低质量的经验数据，进一步无法学习好
- Trust Region一类方法限制一步更新前后策略的差距（用KL散度），进而对策略价值做稳步地提升
- PPO在TRPO的基础上进一步通过限制importance ratio的range，构建优化目标的下界，进一步保证优化的稳定效果，是目前最常用的深度策略梯度算法

THANK YOU