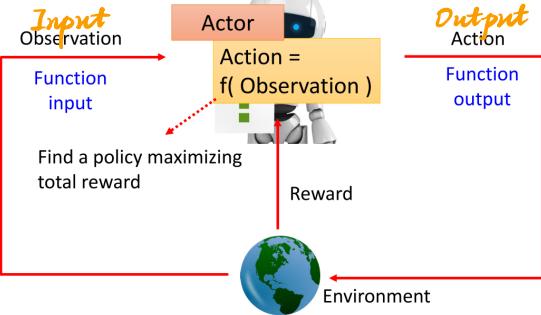


What is RL?



目标：找到一个最大化 Reward 的 Actor
围棋只有赢了才有 Reward=1，输了 Reward=-1，中间大部分 Action 的 reward=0



Step 1: Function with Unknown

→ RL 中的 Actor 通常为 Policy Network

- Input: Observation (向量化表示为 Vector/Matrix)
- Output: 每个 Action 所对应的分值
- 抽样从 Action 中 Sample 出行为 (在面临相同)
- 有时也会采用 Softmax (情况时作出不同动作)

Step 2: Define "Loss"

- 一局游戏结束 = 1 episode (满足结束条件)
- Total Reward ⇔ Return
1' Action 获得的奖赏
- Loss = -Total Reward 一局游戏获得的总 reward

$$\text{Total reward (return): } R = \sum_{t=1}^T r_t$$

What we want to maximize

Step 3: Optimization

- Trajectory $\tau = \{s_1, a_1, s_2, a_2, \dots\}$
- Reward Function 需要同时考虑 s 和 a
- Actor 的 Output 具有随机性
同样的 s 每次产生的 a 不一定相同
- Environment 和 Reward 都是 Black Box 且具有随机性 不是 Network, 无法采用 Gradient Descent
调整网络参数

Random Seed 的随机性：

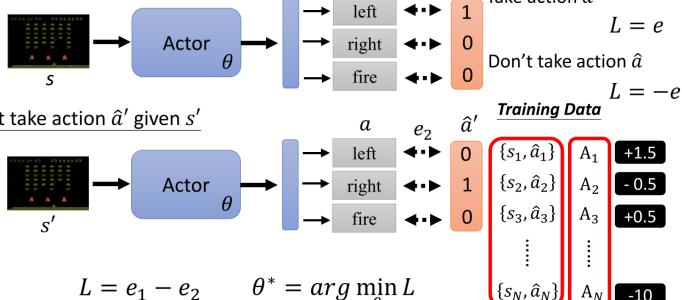
训练时采用不同的初始化参数

RL 的随机性：即使给定相同的 Input, 同一个 agent 的 Output 也可能不同

Policy Gradient

How to Control the Actor

Take action \hat{a} given s



Version 1 从当时时刻往后的

Training Data

$s_1, s_2, s_3, \dots, s_N$

$a_1, a_2, a_3, \dots, a_N$

$r_1, r_2, r_3, \dots, r_N$

所有 reward 之和

$G_1 = r_1 + r_2 + r_3 + \dots + r_N$

$G_2 = r_2 + r_3 + \dots + r_N$

$G_3 = r_3 + \dots + r_N$

cumulated reward

$$G_t = \sum_{n=t}^N r_n$$

距离较远的 action 影响比较小

$$G'_t = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{N-t} r_N$$

Discount factor $\gamma < 1$

越早的 action 累积越高 score

Discounted cumulated Reward

Version 2 动作距离越远，衰减越多

Training Data

$s_1, s_2, s_3, \dots, s_N$

$a_1, a_2, a_3, \dots, a_N$

$r_1, r_2, r_3, \dots, r_N$

Also the credit of a_1 ?

$G_1 = r_1 + r_2 + r_3 + \dots + r_N$

$G'_1 = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$

Discount factor $\gamma < 1$

$$G'_t = \sum_{n=t}^N \gamma^{N-t} r_n$$

How to do the optimization here is the main challenge in RL.

c.f. GAN

$$R(\tau) = \sum_{t=1}^T r_t$$

Version 0 从每个 Action 的 Reward 作 A Training Data → 短视

action 会忽略后续的 observation 及其 reward

many episodes

Short-sighted Version!

Reward Delay

Version 3 标准化 Reward Training Data

设定 baseline 或 G'_t 有正负

Baseline b 平均实力

$A_1 = G'_1 - b$

$A_2 = G'_2 - b$

$A_3 = G'_3 - b$

Good or bad reward is "relative"

If all the $r_n \geq 10$

$r_n = 10$ is negative ...

Minus by a baseline b ???

$$G'_t = \sum_{n=t}^N \gamma^{N-t} r_n$$

Make G'_t have positive and negative values

Policy Gradient

• Initialize actor network parameters θ^0

• For training iteration $i = 1$ to T

- Using actor θ^{i-1} to interact
- Obtain data $\{s_1, a_1\}, \{s_2, a_2\}, \dots, \{s_N, a_N\}$
- Compute A_1, A_2, \dots, A_N
- Compute loss L
- $\theta^i \leftarrow \theta^{i-1} - \eta \nabla L$

Data collection is in the "for loop" of training iterations.

Trajectory of θ^{i-1}

s_1	s_2	s_3	\dots	s_N
a_1	a_2	a_3	\dots	a_N
r_1	r_2	r_3	\dots	r_N

To train θ^i

Exploration Collection Training Data

Suppose your actor always takes "left".

We never know what would happen if taking "fire".

Data Collection 时需要有一些随机性

否则无法得到 Action 可能获得的 Reward

→ Sample Actions 主要原因

- 一般训练是先收集 Data 再进入循环

R2 是跑一次循环 收集一次 Data

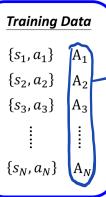
→ 训练过程十分耗时

第一次的资料是旧的经验，对于目前的训练不一定要帮助 (被训练)

→ 收集资料的 Actor 最好和 执行的相同 (和 Environment 互动) → On-policy

不相同 → Off-policy 不用每次重新收集

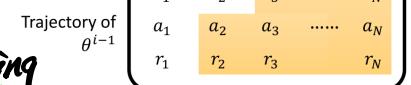
→ 用旧的 Trajectory 训练自己 Training Data



only update once
只更新一次参数

May not observe by theta^i

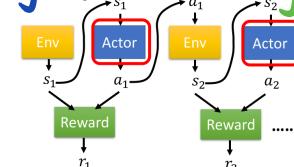
可能由 theta^i 观察不到



但被训练的 Actor 要和自己和示范的不同

How to Exploration:

扩大 Output Entropy
向参数中加入 Noise



Enlarge output entropy
Add noises onto parameters



Critic's
Output 取决于正在评估的 actor

Actor-Critic

Critic Actor θ 观察 State s 和采取 action a 的好坏 Value Function $V^\theta(s)$ $G'_1 = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$

使用 actor θ 时，在 state s 以后能获得的 G' 的期望

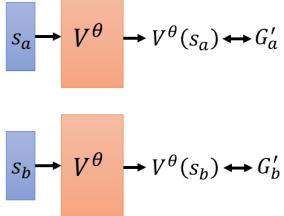
How to Estimate $V^\theta(s)$ -> Critic's Value Function

→ Monte-Carlo (MC) based approach 蒙特卡洛法 → Temporal-difference (TD) Approach 时序差分法 不用等到整个 Episode 结束

Critic 通过 actor θ 和 Environment 互动

After seeing s_a ,

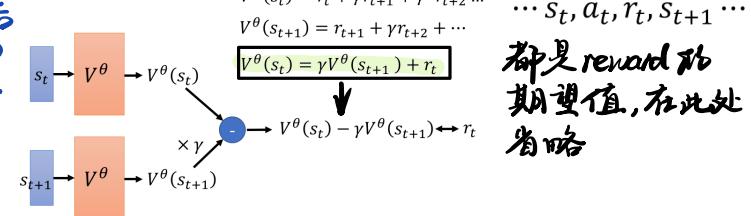
Until the end of the episode, the cumulated reward is G'_a



一个 Episode 结束后所得的 DCR G' 作为训练的 GT

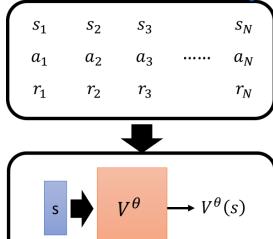
After seeing s_b ,

Until the end of the episode, the cumulated reward is G'_b



都是 reward 和期望值，在此省略

Version 3.5 用 Value Function 评估的基于 (s_t, a_t) 的 G' 的期望作为 Version 3 的 Baseline



Training Data

s_1	s_2	s_3	\dots	s_N
a_1	a_2	a_3	\dots	a_N
r_1	r_2	r_3	\dots	r_N

$\{s_1, a_1\}$	$A_1 = G'_1 - V^\theta(s_1)$
$\{s_2, a_2\}$	$A_2 = G'_2 - V^\theta(s_2)$
$\{s_3, a_3\}$	$A_3 = G'_3 - V^\theta(s_3)$
\vdots	\vdots
$\{s_N, a_N\}$	$A_N = G'_N - V^\theta(s_N)$

$$A_t = G_t - V^\theta(s_t)$$

$$G_t = r_t + \gamma V^\theta(s_{t+1})$$

$$A_t = G_t - V^\theta(s_t) = r_t + \gamma V^\theta(s_{t+1}) - V^\theta(s_t)$$

$$\rightarrow \text{Advantage Actor-Critic (AAC)}$$

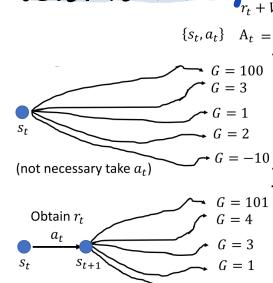
直接用 Critic 决定采用哪个 action

→ DQN <https://arxiv.org/abs/1710.02298>

Rainbow

Actor 和 Critic 部分参数可以共享

Version 4.0 平均减平均代替 sample 成平均



$$\{s_t, a_t\} \quad A_t = G'_t - V^\theta(s_t)$$

→ 在 s_t 下执行完 a_t 以后，直到 Episode 结束，期望获得的 G' 的平均值

$$G_t' = r_t + \gamma V^\theta(s_{t+1})$$

$$A_t = G_t' - V^\theta(s_t) = r_t + \gamma V^\theta(s_{t+1}) - V^\theta(s_t)$$

→ Advantage Actor-Critic (AAC)

直接用 Critic 决定采用哪个 action

→ DQN <https://arxiv.org/abs/1710.02298>

Rainbow

Actor 和 Critic 部分参数可以共享

Reward Shaping

Sparse Reward

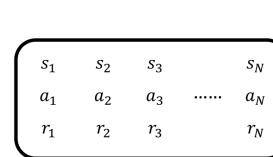
reward 在多数情况下都是 0，每个 action 的执行结果区别不大 (e.g. 机械手臂拧螺丝)

$$A_t = r_t + V^\theta(s_{t+1}) - V^\theta(s_t)$$

→ 是奖励之外的 reward 帮助 agent 学习 Reward Shaping

Curiosity 当 agent 发现新的 (且有意义) 的事物则给予额外 reward

<https://arxiv.org/abs/1705.05363>



Training Data

s_1	s_2	s_3	\dots	s_N
a_1	a_2	a_3	\dots	a_N
r_1	r_2	r_3	\dots	r_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1
$\{s_2, a_2\}$	A_2
$\{s_3, a_3\}$	A_3
\vdots	\vdots
$\{s_N, a_N\}$	A_N

$\{s_1, a_1\}$	A_1

<tbl_r cells="2" ix="3" maxcspan="1" maxrspan="1

No Reward : Learning from Demonstration

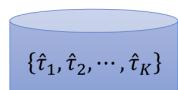
Motivation

- 在真实环境中 reward 的定义和解都非常困难
- 人工定义的 reward 可能导致不可控的行为

Imitation Learning Actor 可以和 Environment 互动，但没有 reward function

→ 使用 Demonstration of the expert 作为参考 (人类的行为方式) → Behavior Cloning (Supervised Learning)

We have demonstration of the expert.



Each $\hat{\tau}$ is a trajectory of the expert.

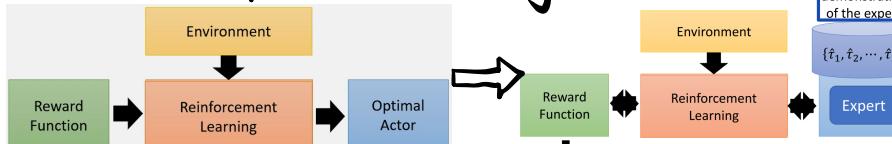
Self driving: record human drivers

Robot: grab the arm of robot

Problem:

- Experts 只会 sample 有限的 Observation (不会做过简单的)
- Agent 会学习一些无关的操作

Inverse Reinforcement Learning 机器自己定义 Reward



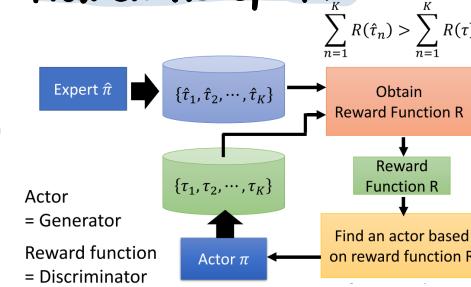
先定义 Reward Function 找到 Optimal Actor

Algorithm

Principle: *The teacher is always the best.*

- Basic idea:
 - Initialize an actor
 - In each iteration
 - The actor interacts with the environments to obtain some trajectories.
 - Define a reward function, which makes the trajectories of the teacher better than the actor.
 - The actor learns to maximize the reward based on the new reward function.
- Output the reward function and the actor learned from the reward function

Framework of IRL



从 Expert Demonstration 反推 Reward Function
 让模仿 Actor 使用 Reward Function 找到 Optimal Actor 可以简单，也可能寻出比较复杂的 Actor

GAN

