

综述

Multi-Agent Reinforcement Learning

MARL 算法可分为 fully cooperative \rightarrow 合作优化总收益

fully competitive \rightarrow 通常收益和为0

和 a mix of the two \rightarrow 两者都有

MARL 的三个特点...

who knows what.

① Multi-dimensional

多维度且对齐，难以评估即定义收益

② 每个 agent 只看重自己的收益 \rightarrow Non Stationary

③ 联合 action space β 通过 agent 数量指数级增加

Background

2.1 Single-Agent RL — 常由 MDP 建模

定义 ① MDP 定义为 (S, A, P, R, γ)

\downarrow \downarrow \downarrow \downarrow
state action $S \times A \rightarrow S$ reward discount factor
transition probability $0 > 1$

$s_t \xrightarrow{a_t} s_{t+1} \sim P_c(s_t, a_t) \rightarrow$ agent 得到 reward
为 $R(s_t, a_t, s_{t+1})$

Core Goals: 寻找 policy $\pi: S \rightarrow \Delta(A)$ 即 S 到 A 分布的 mapping
 $a_t \sim \pi(\cdot | s_t)$

方法: 最大化 $E \left(\sum_{t \geq 0} \gamma^t R(s_t, a_t, s_{t+1}) \mid a_t \sim \pi(\cdot | s_t), s_0 \right)$

State-action function / Q function

$$Q_\pi(s, a) = E \left(\sum_{t \geq 0} \gamma^t R(s_t, a_t, s_{t+1}) \mid a_t \sim \pi(\cdot | s_t), a_0 = a \right)$$

$$V_\pi(s) = E \left(\sum_{t \geq 0} \gamma^t R(s_t, a_t, s_{t+1}) \mid a_t \sim \pi(\cdot | s_t), \begin{matrix} s_0 = s \\ s_0 = s \end{matrix} \right)$$

可由 BP 和 DP 找得 optimal policy.

分类 (1) Value - Based Methods

为寻找最好的 Q-function

policy 则是由 Q-function 经过 greedy 得到

update 公式

$$\hat{Q}(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha [r + \gamma \max_{a'} \hat{Q}(s', a')]$$

on-policy 解决策略与学习率
 \uparrow space learning rate reward
 SARSA, DQN 只是 target 不同 DQN 用 NN 估 a'

另一种 algorithm 为 MCTS 蒙特卡洛树搜索, UCB

还有一种为用已知 Policy 估 value function, TD 算法
 policy evaluation 时序差分

(2) Policy Based Methods 直接在 Policy Space

用差分化函数如 NN 估计 π 才能求

$\pi(\cdot|s) \approx \pi_\theta(\cdot|s)$ 用梯度更新行为 Policy Gradient

PG 的 closed-form \rightarrow

$$\nabla J(\theta) = \mathbb{E}_{a \sim \pi_\theta(\cdot|s), s \sim \eta_{\pi_\theta}(\cdot)} [Q_{\pi_\theta}(s, a)^T \log \pi_\theta(a|s)]$$

J_θ 和 Q_{π_θ} 是 reward 和在 Policy π_θ 下的 Q function 的期望

PG 的常用算法 REINFORCE, G(PO)MDP, actor-critic

非 gradient 方法, 有用 PPO, TRPO 和 soft actor-critic

总的来说 gradient based 比 value based 要好收敛

2.2 Multi-Agent RL Framework

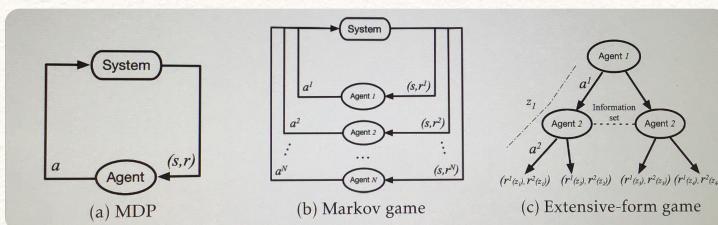
系统中每个 agent 的状态 state 和 reward 的更新都受其他所有 agent 的联合动作

每个 agent 有自己要优化的长期收益, 为其他所有 agents 的延缓

可以将 MARL 分为两类 frameworks

{ Markov/stochastic games
extensive-form games

3个对比



① Markov/Stochastic Games action 完全信息静态博弈。
此时所有 agents 同时选择 a^i after 观察状态 s , 得到独立的 individual reward r^i

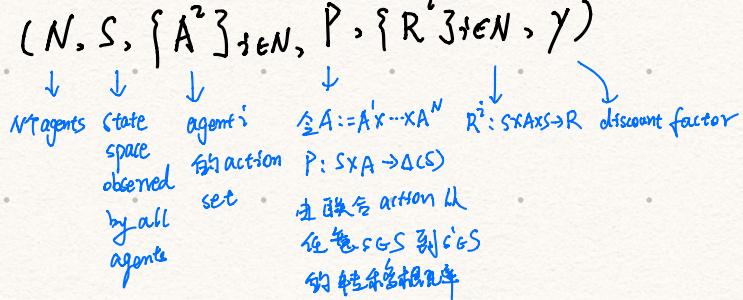
② Extensive-form games 可用于非完全信息博弈。

(此时 agents 选择 a^i 在 game 结束收到 $r^i(z)$)

→ 完全/非完全静态博弈

① MG 的部分定义

Markov Games 定义为一个 tuple $(N, S, \{A^i\}_{i \in N}, P, \{R^i\}_{i \in N}, \gamma)$



At time t , each agent $i \in \mathcal{N}$ executes an action a_t^i , according to the system state s_t . The system then transitions to state s_{t+1} , and rewards each agent i by $R^i(s_t, a_t, s_{t+1})$. The goal of agent i is to optimize its own long-term reward, by finding the policy $\pi^i: \mathcal{S} \rightarrow \Delta(\mathcal{A}^i)$ such that $a_t^i \sim \pi^i(\cdot | s_t)$. As a consequence, the value-function $V^i: \mathcal{S} \rightarrow \mathbb{R}$ of agent i becomes a function of the joint policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ defined as $\pi(a|s) := \prod_{i \in \mathcal{N}} \pi^i(a^i|s)$. In particular, for any joint policy π and state $s \in \mathcal{S}$,

$$V_{\pi^i, \pi^{-i}}^i(s) := \mathbb{E} \left[\sum_{t \geq 0} \gamma^t R^i(s_t, a_t, s_{t+1}) \mid a_t^i \sim \pi^i(\cdot | s_t), s_0 = s \right], \quad (2.3)$$

where $-i$ represents the indices of all agents in \mathcal{N} except agent i . Hence, the solution concept of MG deviates from that of MDP, since the *optimal* performance of each agent is controlled not only by its own policy, but also the choices of all other players of the game.

7

常见解的形式 Nash equilibrium (NE) 纳什平衡的定义如下

定义: MG 的纳什平衡即求得最优联合策略 $\pi^* = (\pi^{1*}, \dots, \pi^{N*})$

$$V_{\pi^{1*}, \pi^{2*}}^1(s) \geq V_{\pi^1, \pi^{2*}}(s) \text{ for any } \pi^1$$

基于 MG 的九种方法:

(1) Cooperative Setting (此时所有 agent 有共同的 reward $R^1 = R^2 = \dots = R^N = R$)

→ 多代理 MDPs in AI 领域
Multiagent MDPs

Markov teams / teams Markov games in Game theory

纳什均衡高

此设置中和 Value function 对每个 agent 都相同, 可以用 single agent 的 optimal policy
General - 是每个 agent 的 reward 不同, 称为 team-average reward (P) 每个 agent 有自己的
reward function 但最终 goal 为优化 long-term average reward

使 decentralized MARL 第三方应用 $\bar{R}(s, a, s') := N^{-1} \sum_{i \in N} R^i(s, a, s')$ for any (s, a, s')

这使得 MARL 的 communication 和 communication efficient 也很必要。

(2) Competitive Setting.

在 MG 中完全 competitive 的建模为 zero-sum Markov games

即 $\sum_{s' \in S} R^i(s, a, s') = 0$ for any (s, a, s')

常考虑 2 个 agents 竞争, 一个 agent 的 reward 就是另一个的 loss

Goal: 优化最差的长期收益

(3) Mixed setting General sum game setting Open AI

此时 goal 和 agents 间的联系和激励有限制。

每个 agent 只关心自己的收益

2.2 Extensive-Form Games

Markov Games 无法处理非完美信息情况，所以用 Extensive-form game