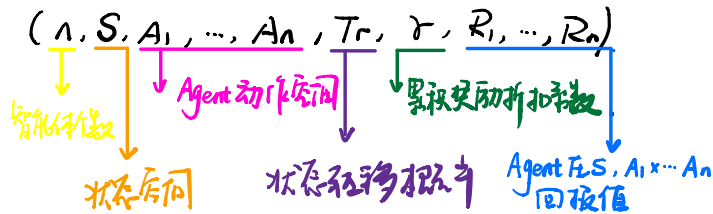


MARL 随机博弈过程描述:

分解为 Stage Game 策略组合

与环境交互更新每个状态阶段博弈中 Q 值

马尔可夫决策过程:



$$Tr: S \times A_1 \times \dots \times A_n \times S' \rightarrow [0, 1]$$

$$R_i: S \times A_1 \times \dots \times A_n \times S' \rightarrow \mathbb{R}$$

MARL 目标:

在每个状态下 纳什均衡策略 (最优策略)。



$$V_i(\pi_1^*, \pi_2^*, \dots, \pi_i^*, \dots, \pi_n^*) \geq V_i(\pi_1^*, \pi_2^*, \dots, \pi_i, \dots, \pi_n^*)$$

$$\forall \pi_i \in \Pi_i \quad i=1, \dots, n$$

所有智能体的联合策略

在纳什均衡处, 对于 All Agent 都不能改变自己策略

从而获得更大的奖励。

最优策略数学描述:

$$(\pi_1^*, \pi_2^*, \dots, \pi_i^*, \dots, \pi_n^*) \quad i=1, \dots, n$$

对于 $\forall s \in S$ 满足

① 基于状态价值函数 $V_i(s)$

$$V_i(\pi_1^*, \pi_2^*, \dots, \pi_i^*, \dots, \pi_n^*) \geq$$

$$V_i(\pi_1^*, \pi_2^*, \dots, \pi_i, \dots, \pi_n^*) \quad \forall \pi_i \in \Pi_i \quad i=1, \dots, n$$

② 基于状态-动作价值函数 $Q_i(s, a_1, \dots, a_n)$

$$\sum_{a_1, \dots, a_n \in A_1 \times \dots \times A_n} Q_i^*(s, a_1, \dots, a_n) \pi_1^*(s, a_1) \dots \pi_i^*(s, a_i) \dots \pi_n^*(s, a_n) \geq$$

$$\sum_{a_1, \dots, a_n \in A_1 \times \dots \times A_n} Q_i^*(s, a_1, \dots, a_n) \pi_1^*(s, a_1) \dots \pi_i(s, a_i) \dots \pi_n^*(s, a_n)$$

Bellman 公式

$$Q_i^*(s, a_1, \dots, a_n) = \sum_{s' \in S} T_i(s, a_1, \dots, a_n, s') [R_i(s, a_1, \dots, a_n, s') + \gamma V_i^*(s')]]$$

$$V_i^*(s) = \sum_{a_1, \dots, a_n \in A_1 \times \dots \times A_n} Q_i^*(s, a_1, \dots, a_n) \pi_1^*(s, a_1) \dots \pi_n^*(s, a_n)$$

如何得到最优策略呢?

IMARL 算法原理

MARL 随机博弈过程
从 Reward 角度进一步分类:

针对不同博弈过程,
设计 MARL 回报函数:

- | | |
|--|---|
| ① fully cooperative ($R_1 = \dots = R_N$) | ① 共同目标, 可以共同最大化. |
| ② fully competitive ($n=2 \quad R_1 = -R_2$) | } \Rightarrow i) 自身学习动态行为的 <u>稳定性</u>
ii) 对其他代理动态行为的 <u>适应</u> |
| ③ Mixed | |

{ Stability: 趋于平稳的策略。 必然性
Adaptation: 其他 Agent 改变策略时 保持或改善性能。 合规性



自然地, 最优策略, Nash equilibrium

收敛收敛到纳什均衡与动态SG的性能之间的联系尚不明确。

聚焦稳定性的算法, 一般认为 Agents 独立的。

适应能力算法, 考虑其它 Agents 的行为。

只考虑稳定性, 不考虑收敛性, 对其他智能体的跟踪。

算法.

Task type → ↓ Agent awareness	Cooperative	Competitive	Mixed
Independent	coordination-free	opponent-independent	agent-independent
Tracking	coordination-based	—	agent-tracking
Aware	indirect coordination	opponent-aware	agent-aware

事件类型, 对其它 Agent 感知.

Task type		Cooperative	Competitive	Mixed
Agent awareness	Independent			
Agent 之间相互独立				
考虑其它 Agent 行为				