

# 「智能博弈对抗方法」最新2022综述-博弈论与强化学习综合视角对比分析

专知  
www.zhuanzhi.ai 专业可信的AI知识分发服务

8人赞同了该文章

## 智能博弈对抗方法:博弈论与强化学习综合视角对比分析

袁唯淋 罗俊仁 陆丽娜 陈佳星 张万鹏 陈璟  
国防科技大学智能科学学院 长沙 410073  
(yuanweilin12@nudt.edu.cn)

智能博弈对抗是人工智能认知决策领域亟待解决的前沿热点问题.以反事实后悔最小化算法为代表的博弈论方法和以虚拟自博弈算法为代表的强化学习方法,依托大规模算力支撑,在求解智能博弈策略中脱颖而出,但对两种范式之间的关联缺乏深入发掘.文中针对智能博弈对抗问题,定义智能博弈对抗的内涵与外延,梳理智能博弈对抗的发展历程,总结其中的关键挑战.从博弈论和强化学习两种视角出发,介绍智能博弈对抗模型、算法.多角度对比分析博弈理论和强化学习的优势与局限,归纳总结博弈理论与强化学习统一视角下的智能博弈对抗方法和策略求解框架,旨在为两种范式的结合提供方向,推动智能博弈技术前向发展,为迈向通用人工智能蓄力.

[jsj.kx.com/CN/10.11896/jsj...](https://jsj.kx.com/CN/10.11896/jsj.20220101)

### 1. 导论

对抗是人类社会发展与演进的主旋律,广泛存在于人与自然、人与人、人与机器之间,是人类思维活动特别是人类智能的重要体现.人工智能浪潮中,对抗的形式不断发生变化,贯穿计算智能、感知智能和认知智能3个阶段[1].以对抗关系为主的博弈研究,为探索认知智能的关键技术原理提供了有效工具.在认知智能层面,信息环境复杂、对抗对手复杂、策略求解复杂等愈发逼近真实世界的复杂场景应用需求,推动了博弈对抗能力的不断提升.高度保留真实世界特性(巨复杂、高动态、强对抗)的智能博弈对抗技术逐渐成为了金融、经济、交通等民用领域的技术引擎和军事智能化实现的重要助推力.在民用领域,尤其是在保护各种关键公共基础设施和目标的挑战性任务[3]中,智能博弈对抗技术不可或缺,例如博物馆、港口、机场等安全机构部署有限的防护资源,在入口处或者外围路网设置安检口进行警力的巡逻防控[4].在军事领域,智能博弈技术积极推动了指挥与控制的智能化发展[5],美国先后启动了“深绿”[6]、指挥官虚拟参谋[7]、“终身学习机器”“指南针”(COMPASS)等项目,旨在缩短“观察G判断G决策G行动”(OODA)的循环时间.

近年来,在人机对抗场景中,AlphaGo[8]、AlphaStar[9]、Pluribus[10]、Suphx[11]、绝悟[12]等一大批高水平AI在游戏验证平台中战胜了人类玩家,智能博弈发展取得了显著突破.智能博弈技术的巨大成功主要依赖于博弈论和强化学习两种范式的结合[13]:博弈论提供了有效的解概念来描述多智能体系统的学习结果,但主要是在理论上发展,应用于实际问题的范围较窄;深度强化学习算法为智能体的训练提供了可收敛性学习算法,可以在序列决策过程中达到稳定和理性的均衡[14][15].一方面,反事实后悔最小化算法(CounterFactual Regret Minimization, CFR)[16]是一种迭代搜索算法,其依托大规模算力支撑,在求解大规模不完备信息博弈策略中脱颖而出,逐渐成为了智能博弈中博弈论范式下的先进代表性算法之一.另一方面,虚拟自博弈算法(Fictitious Self-Play, FSP)[17]依托大规模分布式计算框架,在求解多智能体系统问题中成为了一种通用的强化学习框架,先后被成功应用到雷神之锤III[18]、星际争霸[19]、王者荣耀[12]、德州扑克[20]等复杂大规模对抗场景.CFR与FSP是博弈范式和强化学习范式中的典型方法,也是连接两种范式的算法基础.本文将深挖博弈范式与强化学习范式的深层联系,为两种范式的结合提供方向,推动智能博弈技术前向发展,为迈向通用人工智能蓄力.

本文第2节简要介绍了智能博弈对抗,包括智能博弈对抗的内涵与外延、智能博弈对抗发展历史以及智能博弈对抗中的关键挑战;第3节介绍了智能博弈对抗模型,包括博弈论的基础模型——扩展式博弈模型和强化学习的基础模型——部分可观随机博弈模型,以及结合扩展式博弈模型与马尔可夫模型的通用模型——观察信息可分解的随机博弈模型,从模型上梳理了博弈理论和强化学习的内

在联系;第4节进行了博弈论与强化学习的对比分析,首先详细梳理了博弈论和 强化学习视角下的典型对抗方法,分别以 CFR 和 FSP 为代 表介绍其具体原理,分析变体改进思路,然后多角度对比分析 博弈理论与强化学习的优缺点,探讨后悔值与值函数等基础 概念的联系,归纳总结博弈理论与强化学习的结合方法和结 合框架;第5节介绍了智能博弈对抗研究前沿,归纳了当前热 点前沿智能博弈问题,分别从面向复杂博弈场景的智能博弈 模型、多智能体场景下博弈论与强化学习融合的智能博弈方 法、结合对手建模的 智能 博 弈 方法,以及结合元学习 的多任务场景泛化4 个角度讨论了智能 博弈 前 沿 研究;最后总结全文。

2. 智能博弈对抗简介

智能博弈对抗的内涵与外延

广义上的智能概念涵盖了人工智能、机 器 智 能、混 合 智 能和群体智能。本文的 智 能 概 念 特 指 认 知 智 能 中 机 器 的 自主决策能力,即机器智能,表现为机器模 拟 人 类 的 行 为、思考方式,通过摄像头、话筒等传感器接收 外 界 数 据,与 存 储 器 中 的 数 据 进 行 对 比、识别,从而进行判断、分 析、推 理、决 策。机器智能智能水平的高低可 分 为 若 干 层 次,如 从 最 简 单 的 应 激 反 射 算 法 到 较 为 基 础 的 控 制 模 式 生 成 算 法,再 到 复 杂 神 经 网 络 和 深 度 学 习 算 法。博弈 对 抗 指 代 以 对 抗 关 系 为 主 的 博 弈,在 冲 突 为 主 的 背 景 下 博 弈 方 (拥 有 理 性 思 维 的 个 体 或 群 体)选 择 行 为 或 策 略 加 以 实 施,并 从 中 取 得 各 自 相 应 的 结 果 或 收 益[2 1]。博 弈 与 对 抗 是 人 类 演 化 进 程 中 的 重 要 交 互 活 动,是 人 类 智 能 和 人 类 思 维 方 式 的 重 要 体 现。这种 交 互 活 动 广 泛 存 在 于 个 体 与 个 体、个 体 与 群 体、群 体 与 群 体 之 间。

智能博弈对抗发展历史



图1 智能博弈对抗发展历程与典型应用

博弈对抗不断推动着智能水平的发展,对抗场景从早期的“图灵测试”到目前的“通用场景”探索,不断向真实世界场 景靠拢。2 0 1 6 年,DeepMind基于深度强化学习和蒙特卡洛树搜 索开发的智能 围棋博弈程序 AlphaGo [ 8 ],以 4 : 1 的分数战胜 了人类顶级围棋选手李世石,这标志着人工智能的发展重点 逐渐由感知智能向认知智能过渡。同年,辛辛那提大学基于 遗传模糊树构建的 AlphaAI空战 系统[2 2]在空战对抗中击败 人类飞行员,这成为了无人系统博弈对抗能力生成的推动性 进展。2 0 1 7 年,DeepMind提出的基于自博弈强化学习的棋 类 AIAlphaZero [ 2 3]可以从零开始自学围棋、国际象棋和将 棋,并击败了 AlphaGo。以围棋为代表的完全信息博弈已基 本得到解决,智能博弈的研究开始转向德州扑克和星际争霸 等不完全信息博弈。同年,阿尔伯塔大学和卡内基梅隆大学 先后开发了智能 德州 扑 克 博 弈 程 序 DeepStack [ 2 4]和 LibraG tus [ 2 5],在人 机 对 抗 中 击 败 了 职 业 玩 家。2 0 1 8 年,DeepG Mind在雷神之 锤III夺 旗 游戏中提出了一种基于种群训 练的多 智能 体 强 化 学 习 框 架[1 8],训练 构 建 的 AIFTW 的 性 能 超 越 了 人 类 玩 家 水 平。随 后,智 能 博 弈 朝 着 多 智 能 体 参 与、通 用 场 景 扩 展 迁 移 等 方 向 不 断 发 展,高 效 海 量 数 据 的 实 时 采 样 (数 据)、大 规 模 算 力 加 速 采 样 和 优 化 (算 力)、大 规 模 集 群 架 构 算 法 (算 法)成 为 了 多 智 能 体 强 化 学 习 成 功 的 关 键。博 弈 均 衡 的 方 法 在 多 智 能 体 博 弈 中 仍 存 在 理 论 上 的 局 限 性,但 基 于 两

人框架的多人 博弈扩展依旧在实验中具有较好表现,如 2 0 1 9 年卡内基梅隆大学的六人德州扑克智能博弈程序 Pluribus [ 1 0 ]击败了多名职业玩家.随后,智能博弈的研究趋势开始形成“高质量对抗数据引导”+“分布式强化学习训练”的模式(如麻将 AISuphx,星际争霸 AIAlphaStar [ 1 9 ],谷歌足球 AI觉悟 GWeKick),并逐渐摆脱先验知识,直接完成“端到端”的学习(如捉迷藏 AI [ 2 6 ]、斗地主 AI DouZero [ 2 7 ]、两人德州扑克 AI AIG phaHoldem [ 2 0 ]). 2 0 2 1 年,DARPA 举办的 AlphaDogFight 挑战赛[ 2 8 ]推动了无人系统博弈对抗能力的提升.另一方面,DARPA 开始布局通用 AI 的探索性项目,推动智能博弈向强人工智能迈进.智能博弈对抗发展历程与典型应用总结如图 1 所示.

智能博弈对抗中的关键挑战

复杂博弈环境难评估

(1)不完全信息与不确定性

环境中的不完全信息与不确定性因素提高了博弈决策的难度.战争迷雾造成的不完全信息问题中,关于其他智能体的任何关键信息(如偏好、类型、数量等)的缺失都将直接影响智能体对世界状态的感知,并间接增加态势节点评估的复杂性.不仅如此,考虑不完全信息带来的“欺骗”(如隐真、示假等 [ 2 9 G 3 1 ])行为,将进一步扩展问题的维度.此外,不确定性引入了系统风险,任何前期积累的“优势”都可能因环境中随机因素的负面影响而“落空”.如何综合评估当前态势进行“风险投资”,以获得最大期望回报,成为了研究的另一个难点.另一方面,在策略评估与演化过程中,如何去除不确定因素带来的干扰[ 3 2 ]成为了“准确评价策略的好坏、寻找优化的方向”的难点.

(2)对抗空间大规模

在一些复杂博弈环境中,状态空间和动作空间的规模都非常庞大(见表 1 ),搜索遍历整个对抗空间,无论是在时间约束上还是在存储空间约束上[ 3 3 ]都难以满足要求.

表 1 典型大规模对抗空间博弈场景

Table 1 Typical game scenarios with large adversarial space

博弈场景	状态空间	备注
桥牌	$10^{67}$	—
斗地主	$10^{83}$	—
德州扑克	$10^{160}$	两人无限注
围棋	$10^{170}$	—
王者荣耀	$10^{600}$	1 对 1 场景
兵棋推演	$10^{793}$	城镇居民地想定 <sup>[38]</sup>
星际争霸	$10^{1685}$	128×128 地图,仅考虑 400 个单元的位置

模型抽象[ 3 4 G 3 5 ]的方法在一定程度上可以降低问题的规模,但缺乏理论保证,往往以牺牲解的质量为代价[ 3 6 ].即使以求解次优策略为目标,部分优化算法(如 EGT [ 3 7 ]、一阶 (FirstOrder)算法)仍旧难以直接应用到抽象后的模型.蒙特卡洛采样可以有效地加快算法的速率,但在复杂环境下,如何与其他方法结合并减小搜索中的方差依旧是研究的难点.

多智能体博弈难求解

(1)均衡特性缺失

纳什均衡作为非合作博弈中应用最广泛的解概念,在两人零和场景中具有成熟的理论支撑,但扩展到多智能体博弈时具有较大局限性.两人零和博弈具有纳什均衡存在性和可交换性等一系列优良特性 [ 3 9 ].然而,多人博弈的纳什均衡解存在性缺乏理论保证,且计算复杂,两人一般和博弈的纳什均衡是 PPAD 难问题[ 4 0 ],多人一般和的计算复杂度高于 PPAD.即使可以在多人博弈中有效地计算纳什均衡,但采取这样的纳什均衡策略并不一定是“明智”的.如果博弈中的每个玩家都独立地计

算和采取纳什均衡策略,那么他们的策略组合可能并不是纳什均衡,并且玩家可能具有偏离到不同策略的动机[41G42].

## (2) 多维学习目标

对于单智能体强化学习而言,学习目标是最大化期望奖励,但是在多智能体强化学习中,所有智能体的目标不一定是一致的,学习目标呈现出了多维度[13].学习目标可以分为两类[43]:理性和收敛性.当对手使用固定策略时,理性确保了智能体尽可能采取最佳响应,收敛性保证了学习过程动态收敛到一个针对特定对手的稳定策略,当理性和收敛性同时满足时,会达到新的纳什均衡.

## (3) 环境非平稳

当多个智能体同时根据自己的奖励来改进自身策略时,从每个智能体角度来看,环境变得非平稳,学习过程难以解释[44].智能体本身无法判断状态转移或奖励变化是自身行为产生的结果,还是对手探索产生的.完全忽略其他智能体独立学习,这种方法有时能产生很好的性能,但是本质上违背了单智能体强化学习理论收敛性的平稳性假设[45].这种做法会失去环境的马尔可夫性,并且静态策略下的性能测度也随之改变.例如,多智能体中单智能体强化学习的策略梯度法的收敛结果在简单线性二次型博弈[46](Linear Quadratic Games)中是不收敛的.

# 3 智能博弈对抗模型

## 扩展式博弈模型

扩展式博弈适用于序贯决策中建模智能体与环境的重复交互过程,尤其是存在“智能体对其他智能体之前的决策节点不可分辨(含有隐藏信息)”或者“智能体遗忘之前的决策(不完美回忆)”的情景.

## 部分可观随机博弈

与扩展式博弈的树结构不同,马尔可夫博弈(Markov Game)也称随机博弈(Stochastic Game),具有马尔可夫链式结构.

## 通用模型

博弈理论和强化学习理论并不互斥,在模型上,博弈论的扩展式博弈模型和强化学习的部分可观随机博弈两种模型之间具有一定的联系,例如都可以通过放宽某些条件限制转化为观察信息可分解的随机博弈(Factored Observation Stochastic Games, FOSG)[58].FOSG是POSG的一种扩展性变体,模型聚焦于公共信息(Public Information)的表示和分解,如图3所示.

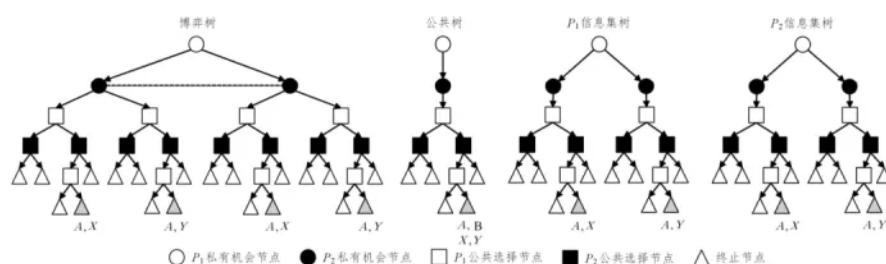


图3 观察信息分解博弈树<sup>[50]</sup>

# 4 博弈论与强化学习的对比分析

博弈论和强化学习是求解序贯决策问题的有效工具,然而它们在算法特性(泛化性、可解释性、收敛性保证)、应用场景(多人博弈、序贯博弈、即时策略博弈)以及硬件资源(算力需求)等方面各有所长,本文总结了近5年AAAI,IJCAI,NeuralPS,AMMAS,ICRL等人工智能顶刊顶会中与智能博弈技术相关的博弈论与强化学习文章,按专家打分的方法,绘制对比分析雷达图,如图7所示.博弈理论在两人零和博弈问题上已经具有较为成熟的理论,包括纳什均衡(以及其他解概念)的等价性、存在性、可交换性(Interchangeability)[39]等,但在多人博弈问题中还需要新的解概念以及相关理论的支持.CFR算法通过后悔值迭代更新生成策略,模型具有可解释性.但是,完美回放和终端可达的



强烈假设限制了 CFR 的使用场景[9 2].强化学习结合深度学习,直接实现端到端的学习,具有很强的泛化性,在多智能体博弈中已取得较多成功应用.但网络的训练往往需要超大规模的算力支撑,且模型的可解释性不强.本节将对两种方法的具体局限性进行深入剖析,为两种方法的结合互补提供方向.

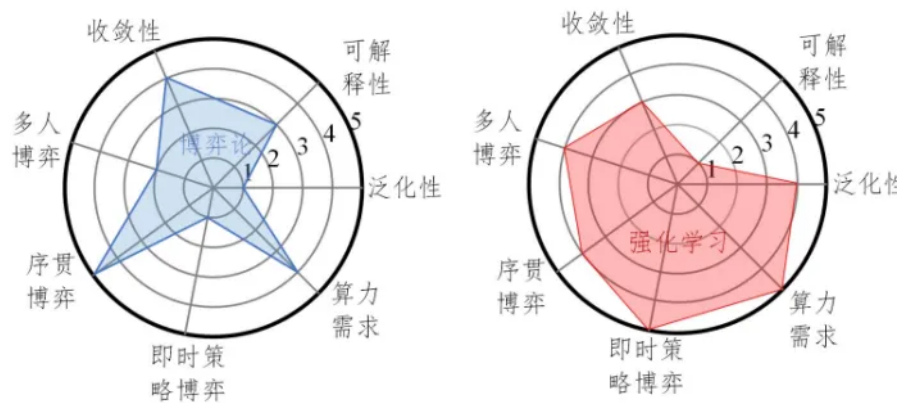


图 7 博弈论与强化学习方法特性对比

## 5. 智能博弈对抗研究前沿

### 面向复杂博弈场景的智能博弈模型

博弈论中的“信息集”和强化学习中的“观察函数”都是针对智能博弈场景中不完全信息的形式化描述.现实世界中,博弈场景更加复杂,不完全信息引发了博弈信息“不可信”等问题——智能体通常是不完全理性的,并且不同认知层次的智能体参与的博弈具有“欺诈[1 1 4 G 1 1 5]”“隐藏”“合谋”“认知嵌套(建模对手的同时,对手也在建模利用己方)”[1 1 6 G 1 1 7]等新挑战.如何针对认知博弈对抗中的新挑战,形式化描述“欺骗”等复杂博弈要素,建立复杂博弈信息的量化评估体系,成为了智能博弈向真实世界迁移应用的模型基础.

### 多智能体场景下博弈论与强化学习融合的智能博弈方法

虽然博弈论提供了易于处理的解决方案概念来描述多智能体系统的学习结果,但是纳什均衡是一个仅基于不动点的静态解概念,在描述多主体系统的动态特性方面(如循环集 (RecurrentSet)[1 1 8]、周期轨道 (Periodic Orbits)和极限环 (LimitCycles)[5 5])具有局限性.一方面,寻求具有更多优良特性的多人博弈新解概念,探索多人局部纳什均衡点求解方法,将是博弈视角下求解多智能体博弈问题的新突破口;另一方面,发挥深度学习和强化学习在信息表征、复杂函数拟合方面的优势,基于自博弈求解框架、值函数评估方法、强化学习结合 CFR等方法,探索博弈论方法与强化学习的有效融合机制,将是突破多智能体博弈学习瓶颈的前沿方向.

### 结合对手建模的智能博弈方法

对抗胜负的本质是超越对手的相对优势,决策的制定必须以对手的行动或策略为前提.纳什均衡是应对未知通用对手时最小化最坏可能性,用最“保险”的策略应对,而并不是寻求最优应对策略.放宽纳什均衡中“未知通用对手”的设定,考虑不完全理性对手的最佳应对,一些新的解概念[6 1]被提出,结合显式对手建模(ExplicitOpponentModeling)[1 1 9]和均衡近似,平衡利用性与剥削性,实现多目标优化,为融合对手建模的博弈学习提供参考.此外,在一些更加复杂的对抗场景中,如对手具有学习意识 (OpponentG Learning AwareG ness)[1 2 0]等,最大熵多智能体强化学习[1 2 1]成为研究如何进行博弈均衡对抗策略选择的新趋势.不仅如此,反对手建模方法也在同步发展.基于意图识别设计[1 2 2]的对抗意图识别、包含意图隐藏与欺骗的对抗意图识别方法等[1 2 3]反对手建模方法在欺骗路径规划[2 9]等问题中得到了一定的研究.在复杂博弈对抗场景中,如何基于对手模型安全利用对手,以及如何保全自我反对手建模成为了新的探索性研究.

结合元学习的多任务场景泛化

学习模型如何更好地泛化到差异很大的新领域中,是一种更加高效和智能的学习方法.元学习(MetaLearning)逐渐发展成为让机器学会学习的重要方法.元学习是通用人工智能(GeneralAI)的分支,通过发现并推广不同任务之间的普适规律来解决未知难题.元学习的输入是一个任务集合,目的是对每个任务的特性和任务集合的共性建模,发现任务之间的共性和内在规律,以追求在差异较大任务之间的迁移,且不会产生较大的精度损失,目前其已经扩展到元强化学习[1 2 4 G 1 2 5]、元模仿学习[1 2 6]、元迁移学习、在线元学习[1 2 7]、无监督元学习[1 2 8 G 1 2 9]等.如何结合博弈理论和元强化学习的优势,构建高效、可解释性强、具有收敛性保障和泛化性的近似纳什均衡求解体系,将是未来智能博弈技术发展的巨大推动力之一.

结束语

本文针对智能博弈对抗问题,介绍了智能博弈对抗的内涵与外延,梳理了智能博弈对抗发展历程,总结了其中的关键挑战.从博弈论和强化学习两种视角出发,介绍了智能博弈对抗模型和算法,多角度对比分析了博弈理论和强化学习的优势与局限,归纳总结了博弈理论与强化学习统一视角下的智能博弈对抗方法和策略求解框架,旨在为两种范式的结合提供方向,推动智能博弈技术前向发展,为迈向通用人工智能蓄力.

智能博弈对抗方法:博弈论与强化学习综合视角对比分析 - 专知VIP

[www.zhuanzhi.ai/vip/1612538fbd96a366a4527f3ce339...](http://www.zhuanzhi.ai/vip/1612538fbd96a366a4527f3ce339...)

「智能博弈对抗方法」最新2022综述-博弈论与强化学习综合视角对比分析

[mp.weixin.qq.com/s/us0AYQg8KRDvawPF...](https://mp.weixin.qq.com/s/us0AYQg8KRDvawPF...)



编辑于 2022-08-29 15:06

博弈论

写下你的评论...

评论内容由作者筛选后展示



还没有评论，发表第一个评论吧

推荐阅读

博弈论(5) 自动机、单步偏离原理、重复博弈

在上一节我们介绍了动态博弈的框架、子博弈精炼纳什均衡与逆向递



归法的关系. 在本节我们讨论一类更具体的动态博弈: 重复博弈. 重复博弈在实证研究中起到了很好的指引作用. 我们看看两方博弈.

TOVAR... 发表于高级微观经...

