

# 基于值分布的多智能体分布式深度强化学习算法<sup>①</sup>



陈妙云, 王 雷, 盛 捷

(中国科学技术大学 信息科学与技术学院, 合肥 230027)

通信作者: 王 雷, E-mail: wangl@ustc.edu.cn

**摘 要:** 近年来深度强化学习在一系列顺序决策问题中取得了巨大的成功, 使其为复杂高维的多智能体系统提供有效优化的决策策略成为可能. 然而在复杂的多智能体场景中, 现有的多智能体深度强化学习算法不仅收敛速度慢, 而且算法的稳定性无法保证. 本文提出了基于值分布的多智能体分布式深度确定性策略梯度算法 (multi-agent distributed distributional deep deterministic policy gradient, MA-D4PG), 将值分布的思想引入到多智能体场景中, 保留预期回报完整的分布信息, 使智能体能够获得更加稳定有效的学习信号; 引入多步回报, 提高算法的稳定性; 引入了分布式数据生成框架将经验数据生成和网络更新解耦, 从而可以充分利用计算资源, 加快算法的收敛. 实验证明, 本文提出的算法在多个连续/离散控制的多智能体场景中均具有更好的稳定性和收敛速度, 并且智能体的决策能力也得到了明显的增强.

**关键词:** 多智能体; 深度强化学习; 值分布; 多步回报; 分布式数据生成

引用格式: 陈妙云, 王雷, 盛捷. 基于值分布的多智能体分布式深度强化学习算法. 计算机系统应用, 2022, 31(1): 145-151. <http://www.c-s-a.org.cn/1003-3254/8237.html>

## Multi-agent Distributed Deep Reinforcement Learning Algorithm Based on Value Distribution

CHEN Miao-Yun, WANG Lei, SHENG Jie

(School of Information Science and Technology, University of Science and Technology of China, Hefei 230027, China)

**Abstract:** In recent years, deep reinforcement learning has achieved great success in many sequential decision-making problems, which makes it possible to provide effective and optimized decision-making strategies for complex and high-dimensional multi-agent systems. However, in complex multi-agent scenarios, the existing multi-agent deep reinforcement learning algorithm has a low continuous convergence speed, and the stability of the algorithm cannot be guaranteed. Herein, we propose a new multi-agent deep reinforcement learning algorithm, which is called multi-agent distributed distributional deep deterministic policy gradient (MA-D4PG). We adapt the idea of value distribution to multi-agent scenarios and retain the complete distribution information of expected return, so that agents can obtain a more stable and effective learning signal. We also introduce a multi-step return to improve the stability of the algorithm. In addition, we use a distributed data generation framework to decouple empirical data generation and network update for the purpose of taking full advantage of computing resources to speed up the convergence. Experiments show that the proposed method has better stability and a higher convergence speed in multiple continuous/discrete controlled multi-agent scenarios and the decision-making ability of agents has also been significantly enhanced.

**Key words:** multi-agent; deep reinforcement learning; value distribution; multi-step return; distributed data generation

<sup>①</sup> 基金项目: 中国科学技术大学预研基金 (YZ2101900004)

收稿时间: 2021-03-11; 修改时间: 2021-04-07, 2021-04-16; 采用时间: 2021-04-20; csa 在线出版时间: 2021-12-17

## 1 引言

近年来,强化学习在解决一系列顺序决策问题上取得了巨大的成功. Minih 等人<sup>[1]</sup>将深度学习和强化学习结合起来,提出了深度强化学习,在众多领域都得到了广泛的应用. 比如机器人控制<sup>[2]</sup>, 生命科学<sup>[3]</sup>, 自动驾驶<sup>[4]</sup>, 5G 通信<sup>[5]</sup>, 物联网<sup>[6]</sup>, 车联网<sup>[7]</sup>等. 在深度强化学习的基础上众多研究者提出了多项改进来提升算法的收敛速度和稳定性. 然而这些算法的研究和应用大多集中在单智能体领域. 现实世界中的任务和场景通常更加复杂, 涉及到多个智能体之间的合作和竞争. 比如工厂中的协作机器人, 交通控制系统和自动军事系统等. 将单智能体强化学习算法简单地移植到多智能体场景中通常会带来一些问题和挑战<sup>[8]</sup>. 比如单智能体深度强化学习算法无法表征智能体之间的合作与竞争关系以及智能体之间共同进化导致环境非平稳性的问题. 并且多智能体与复杂场景会带来计算量大幅提升, 算法收敛速度明显下降等挑战.

目前, Self-play<sup>[9]</sup>被证明是一种有效的多智能体训练范式, 并成功的应用在了围棋, 象棋以及博弈竞争场景中. Palmer 等人<sup>[10]</sup>提出了一种基于宽容的 DQN 算法, 对同伴不合理的行为给予一定的容忍度. Foerster 等人<sup>[11]</sup>提出了反事实的思想, 解决了多智能体合作场景中信用分配的问题. 这些算法在多智能体场景中取得了很好的效果. 但是只能应用在离散动作空间中, 并且 Self-play 还增加了智能体之间动作集完全一致的限制. Lowe 等人<sup>[12]</sup>提出了一种集中式训练, 分散式执行的算法 (multi-agent deep deterministic policy gradient, MA-DDPG), 解决了多智能体场景中环境非平稳的问题. 并且该算法可以同时应用在连续和离散控制的多智能体场景中, 使得应用场景更加灵活.

然而对于复杂的多智能体场景, 目前的多智能体深度强化学习算法还有很多问题没有解决. 智能体动作策略的更新质量非常依赖于预期回报的估计, 但是多个智能体与环境交互时, 其内在的随机性会导致预期回报具有不确定性, 现有的方法将预期回报简单地平均为期望值则会丢失完整的分布信息, 导致智能体不能进行有效的学习. 并且现有的方法主要基于单步回报, 智能体只考虑当下的收益而忽略未来的信息, 从而引发算法不稳定. 此外现有方法中经验数据生成往往和网络更新耦合在一起, 导致只能串行处理, 难以有效利用更多的计算资源. 为此, 本文提出了基于值分布

的多智能体分布式深度确定性策略梯度算法 (multi-agent distributed distributional deep deterministic policy gradient, MA-D4PG), 引入值分布的思想<sup>[13]</sup>对预期回报的分布建模, 避免分布信息丢失, 使得智能体能够获得更加准确的梯度更新, 保证学习的有效性; 引入多步回报使得智能体能够从未来的经验中学习, 增加智能体的预见能力, 提高算法的稳定性; 引入分布式数据生成框架 (APE-X)<sup>[14]</sup>, 将智能体与环境交互生成经验数据部分与网络更新解耦, 使得经验数据生成过程和网络更新能够并行进行, 充分利用更多的计算资源.

在多个不同的连续/离散控制的多智能体场景中本文算法的稳定性和收敛速度都有了极大的提升, 智能体的决策能力也有了明显的增强. 最后本文还做了消融实验展示每个改进各自的贡献以及相互作用.

## 2 相关工作

### 2.1 多智能体马尔可夫决策过程

多智能体马尔可夫决策过程是马尔可夫决策过程在多智能体场景下的扩展. 其同样由四元组构成:

$$M = (S, A, Pr(S'|S, A), R) \quad (1)$$

其中,  $S$ 表示所有智能体的状态集合 $\{s_1, s_2, \dots, s_k\}$ ,  $k$ 表示智能体的数量.  $A$ 表示所有智能体的动作集合 $\{A_1, A_2, \dots, A_k\}$ .  $Pr(S'|S, A)$ 表示状态转移函数,  $S \times A_1 \times \dots \times A_k \rightarrow P(S')$ . 其中 $P(S')$ 定义了所有智能体的下一个状态的概率分布.  $R$ 表示所有智能体的回报集合 $\{r_1, \dots, r_k\}$ . 本文遵循以上符号表示.

对于每一个智能体 $i$ , 学习相应的动作选择策略 $\pi_i: o_i \rightarrow P(A_i)$ . 其中 $o_i$ 表示智能体 $i$ 的局部观测,  $P(A_i)$ 表示选择动作的概率分布.

智能体的目标是学习一个动作选择策略以最大化预期回报如下, 其中 $\gamma$ 为折现因子.

$$J_i(\pi_i) = E_{A_1 \sim \pi_1, \dots, A_k \sim \pi_k} \left[ \sum_{t=0}^{\infty} \gamma^t r_{i,t}(S_t, A_{1,t}, \dots, A_{k,t}) \right] \quad (2)$$

### 2.2 多智能体深度确定性策略梯度

深度策略梯度算法 (deep policy gradient, DPG)<sup>[15]</sup>是一种非常经典的深度强化学习算法. 该算法使用神经网络模拟动作策略函数, 通过梯度上升法调整网络参数以最大化目标函数, 如式 (3) 所示:

$$J(\theta) = Q(S_t, A_t) \quad (3)$$

其中,  $Q(S_t, A_t)$ 表示状态动作对的预期回报的期望值. 相应的动作策略网络的梯度公式如式 (4) 所示:

$$\nabla J(\theta) = E_{S_t, A_t} [\nabla_{\theta} \log \pi(A_t | S_t) Q(S_t, A_t)] \quad (4)$$

其中,  $\pi$  表示智能体的随机动作策略, 输出所有动作的概率分布。

DDPG (deep deterministic policy gradient)<sup>[16]</sup> 是 DPG 的一种变体, 其中动作策略输出的是确定性动作而不是所有动作的概率分布。可以应用于连续控制的多智能体场景中, 相应的动作策略网络梯度公式如式 (5):

$$\nabla J(\theta) = E_{S_t, A_t} [\nabla_{\theta} \mu(A_t | S_t, \theta) \nabla_{A_t} Q(S_t, A_t)] \quad (5)$$

其中,  $\mu$  表示智能体的确定性动作策略。

而多智能体深度确定性策略梯度算法是 DDPG 在多智能体场景的扩展, 并且与 actor-critic 框架相结合。其中 actor 模拟动作选择策略函数  $\mu$ , 根据智能体  $i$  的局部观测  $o_i$ , 输出动作  $A_i = \mu_i(o_i | \theta_i)$ 。critic 模拟动作价值函数  $Q_i$  以输出智能体  $i$  某一时刻的状态动作对的预期回报的期望值。 $Q_i$  是一个集中式的函数, 考虑了所有智能体的状态动作信息, 解决了多智能体场景中智能体共同学习以及相互作用导致的环境非平稳。

在 MA-DDPG 算法智能体  $i$  的 actor 网络的梯度公式以及 critic 网络的损失函数分别如式 (6) 和式 (7):

$$\nabla J_i(\theta_i) = E_{S, A} [\nabla_{\theta_i} \mu_i(A_i | S_i, \theta_i) \nabla_{A_i} Q_i(S, A_1, \dots, A_k)] \quad (6)$$

$$\begin{cases} L_i(w_i) = E_{S, A, r, S'} [(Q_i^r(S, A_1, \dots, A_k | w_i) - y)^2] \\ y = r_i + \gamma \max_{A'} \bar{Q}_i^r(S', A'_1, \dots, A'_k | \bar{w}_i) \end{cases} \quad (7)$$

### 3 基于值分布的多智能体分布式深度确定性策略梯度算法

本文引入了值分布的思想对预期回报的分布建模, 提出了一种基于值分布的评论家 (critic) 更新方法; 并在此基础上引入多步回报衡量 critic 的损失; 最后引入了分布式数据生成架构 (APE-X) 将智能体与环境交互生成经验数据的部分与网络更新解耦。这些将会在本节进行详细描述。

#### 3.1 基于值分布的 critic 更新

基于值的多智能体深度强化学习算法直接对智能体状态动作对的预期回报的期望建模, 不考虑智能体与环境的交互的内在随机性对预期回报的影响。这种建模方式损失了关于预期回报的完整的分布信息。本文引入值分布的思想, 强调预期回报的分布信息的重要性, 对预期回报的分布建模, 而非对其期望值建模。本文将预期回报看作随机变量, 记为  $Z_i^{\mu_i}(S_t, A_t)$ 。定义如

式 (8) 所示:

$$Q_i^{\mu_i}(S_t, A_t) = E_{S_t, A_t, S_{t+1}} [Z_i^{\mu_i}(S_t, A_t)] \quad (8)$$

其中,  $Q_i^{\mu_i}(S_t, A_t)$  表示智能体  $i$  的 critic 网络输出预期回报的期望值。

本文使用参数化的多分类分布对预期回报的分布建模  $\{N \in \mathbb{N}, V_{\min}, V_{\max} \in \mathbb{R}\}$ , 其中  $N$  表示分类个数,  $V_{\min}$  表示最小类别值,  $V_{\max}$  表示最大类别值。以此模拟预期回报的近似分布。其中一系列的类别值分别如式 (9):

$$z_j = V_{\min} + j \cdot \Delta z, \quad 0 \leq j < N, \quad \Delta z = \frac{V_{\max} - V_{\min}}{N - 1} \quad (9)$$

在 actor-critic 的框架下, 改进的基于值分布思想的 critic 网络输出预期回报的  $N$  分类分布  $Z_i^{\mu_i}$ , 表示各类预期回报值的概率分布。记  $p_{i,j}(S_t, A_t)$  为智能体  $i$  的 critic 网络所输出的预期回报属于第  $j$  类  $z_j$  的概率。

$$p_{i,j}(S_t, A_t) = \frac{e^{w_j(S_t, A_t)}}{\sum_{n=0}^{N-1} e^{w_n(S_t, A_t)}}, \quad 0 \leq j < N - 1 \quad (10)$$

其中,  $S_t = \{s_{1,t}, \dots, s_{k,t}\}$  表示所有智能体的状态信息,  $A_t = \{A_{1,t}, \dots, A_{k,t}\}$  表示所有智能体的动作信息。

基于值分布的 critic 的损失函数定义如式 (11) 和式 (12):

$$L_i(w_i) = E[d((T\bar{Z}_i^{\mu_i})(S_t, A_{1,t}, \dots, A_{k,t}), Z_i^{\mu_i}(S_t, A_{1,t}, \dots, A_{k,t}))] \quad (11)$$

$$\begin{aligned} (T\bar{Z}_i^{\mu_i})(S_t, A_{1,t}, \dots, A_{k,t}) = \\ r_{i,t} + \gamma E[Z_i^{\mu_i}(S_{t+1}, A_{1,t+1}, \dots, A_{k,t+1})] \end{aligned} \quad (12)$$

其中,  $Z_i^{\mu_i}$  表示智能体  $i$  的 critic 网络输出的预期回报的近似分布 ( $Z$  分布),  $\bar{Z}_i^{\mu_i}$  表示智能体  $i$  的目标 critic 网络的输出,  $(T\bar{Z}_i^{\mu_i})$  表示值分布的贝尔曼算子。式 (12) 将目标 critic 网络的输出映射成目标  $Z$  分布,  $d$  表示预期  $Z$  分布和目标  $Z$  分布之间距离的度量指标, 本文采用交叉熵损失进行度量。

以上所述引入了值分布的思想对 critic 的更新进行改进, 相应的 actor 更新也引入预期回报的分布信息, 本文通过引入  $Z$  分布的期望来实现, 则 actor 网络更新的梯度公式如式 (13) 所示:

$$\begin{cases} \nabla J_i(\theta_i) = E_{S_t, A_t} [\nabla_{\theta_i} \mu_i(A_{i,t} | o_{i,t}, \theta_i) E[\nabla_{A_{i,t}} Z_i(S_t, A_{1,t}, \dots, A_{k,t})]] \\ A_{i,t} = \mu_i(o_{i,t} | \theta_i) \end{cases} \quad (13)$$

#### 3.2 多步回报

本小节在第 3.1 节所述的基于值分布的 critic 更新

的基础上引入多步回报来计算 critic 的 TD (temporal-difference) 误差. 利用多步回报可以使得智能体从未来的经验中学习, 而不是只考虑当下的收益. 多步回报的思想已经广泛的应用在众多单智能体强化学习算法<sup>[17,18]</sup>中. 类似的修正也可以应用在多智能体深度强化学习算法中. 基于多步回报的贝尔曼方程定义如式 (14) 和式 (15) 所示:

$$(T\bar{Z}_i^{t'}) (S_t, A_{1,t}, \dots, A_{k,t}) = R_{i,t} + \gamma E[\bar{Z}_i^{t'}(S_{t+N}, A_{1,t+N}, \dots, A_{k,t+N})] \quad (14)$$

$$R_{i,t} = \sum_{t'=t}^{t+N-1} \gamma^{t'-t} r_{i,t'} \quad (15)$$

其中,  $R_{i,t}$  表示时刻  $t$  智能体  $i$  的  $N$  步回报.

### 3.3 分布式数据生成和网络更新

针对模型耦合度高无法充分利用计算资源的问题, 本文引入了分布式数据生成的框架 (APE-X) 将多智能体场景中智能体与环境交互生成经验数据和网络更新两部分进行解耦. 负责数据生成的部分称为 Actor, 负责从经验回放池中采样数据进行网络更新的部分称为 Learner. Actor 的网络结构与 Learner 完全一致, 只负责与环境交互得到经验数据存入经验回放池中, 不执行梯度计算进行参数更新, 而是定期复制 Learner 的网络参数. Learner 只负责从经验回放池中采样数据进行网络更新而不负责与环境进行交互. 本文算法并行运行多个独立的 Actor, 分别与环境进行交互得到经验数据  $\{S_{t+1-N}, A_{t+1-N}, R_{t+1-N}, S_{t+1}\}$ , 写入同一个经验缓冲池中. Learner 从该经验缓冲池中随机采样  $M$  批经验数据进行网络更新. 具体的算法框架如图 1 所示. 其中  $S_{t+1-N} = \{s_{1,t+1-N}, \dots, s_{k,t+1-N}\}$ ,  $A_{t+1-N} = \{A_{1,t+1-N}, \dots, A_{k,t+1-N}\}$ ,  $R_{t+1-N} = \{R_{1,t+1-N}, \dots, R_{k,t+1-N}\}$ .

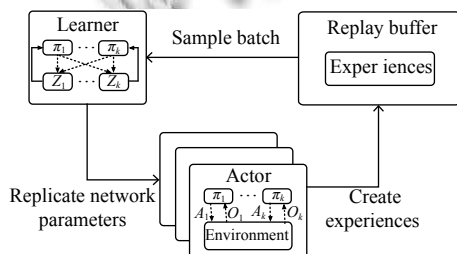


图 1 分布式数据生成框架示意图

### 3.4 本文算法小结

本文算法的伪代码如算法 1 和算法 2 所示. Learner

中 actor-critic 的网络都采用梯度下降法进行网络参数更新. actor 和 critic 的学习率分别是  $\alpha_t, \beta_t$ , 均采用 Adam 优化器<sup>[19]</sup>, 利用动量和自适应学习率来加快收敛速度. 该算法伪代码中同时包含了 Actor 和 Learner 的伪代码, 多个 Actor 初始化为相同的多智能体环境, 并行与环境交互生成经验数据并将其存入经验回放池中.

算法 1. 经验数据生成-Actor

- 1) 初始化  $N$  值, 智能体数量  $K$ , 经验回放池  $B$ , 初始化软更新频率  $TAU$
- 2) 初始化 Actor 网络参数  $(\theta_{Actor}^{1:k}, \omega_{Actor}^{1:k}) \leftarrow (\theta_{Learner}^{1:k}, \omega_{Learner}^{1:k})$
- 3) 初始化 States, Actions, Rewards, Terminals 大小为  $N+1$  的 deque
- 4) **For**  $j=1, \dots, Episode\_num$  **do**
- 5)   Clear (States, Actions, Rewards, Terminals)
- 6)    $S_0 \leftarrow Env.Reset()$
- 7)   States.Add( $S_0$ )
- 8)   **For**  $t=1, \dots, max\_step$  **do**
- 9)      $A_{1,t-1}, \dots, A_{k,t-1} \leftarrow \mu_1(s_{1,t-1}), \dots, \mu_k(s_{k,t-1})$
- 10)     $r_{1,t-1}, \dots, r_{k,t-1}, S_t, term_t \leftarrow Env.Step(A_{t-1} + OU)$
- 11)    States.Add( $S_t$ )  
      Actions.Add( $A_{t-1}$ )  
      Rewards.Add( $\{r_{1,t-1}, \dots, r_{k,t-1}\}$ )  
      Terminals.Add( $term_t$ )
- 12)    **If** States.full **then**
- 13)      根据式 (15) 计算  $N$  步回报  $R_{t-N}$
- 14)      B.Add(States[0], Actions[0],  $R_{t-N}$ , States[N], Terminals[N-1])
- 15)      PopLeft(States, Actions, Rewards, Terminals)
- 16)    **End If**
- 17)    **If**  $term_t$  **break**
- 18)    **End For**
- 19) **End For**

算法 2. 网络更新-Learner

- 1) 初始化采样批次  $M$ , 最大更新步数  $T$
- 2) 随机初始化网络参数  $(\theta_{Learner}^{1:k}, \omega_{Learner}^{1:k})$
- 3) 初始化目标网络参数  $(\theta_{target}^{1:k}, \omega_{target}^{1:k}) \leftarrow (\theta_{Learner}^{1:k}, \omega_{Learner}^{1:k})$
- 4) **For**  $t=1, \dots, T$  **do**
- 5)   **If**  $len(B) < M$  **break**
- 6)    从  $B$  中随机采样  $M$  批数据
- 7)    **For**  $i=1, \dots, k$  **do**
- 8)      根据式 (14) 和式 (11) 计算 critic 的损失, 使用梯度下降法进行反向传播更新 critic 网络参数
- 9)      根据式 (13) 计算 actor 的梯度, 其中  $A_{i,t} = \mu_i(s_{i,t} | \theta_i)$ , 更新 actor 网络参数
- 10)    **End For**
- 11)     $(\theta_{target}^{1:k}, \omega_{target}^{1:k}) = (1-TAU) \cdot (\theta_{target}^{1:k}, \omega_{target}^{1:k}) + TAU \cdot (\theta_{Learner}^{1:k}, \omega_{Learner}^{1:k})$
- 12)    **If**  $t \bmod actor\_update\_freq == 0$  **then**
- 13)      Replicate Network Weights to Actors
- 14)    **End If**
- 15) **End For**

## 4 实验分析

本文设置了3个实验场景来测试本文算法(MA-D4PG)的性能,将在本节依次描述.在实验中设置了一个额外的Actor.该Actor只利用,而不进行探索,方便在没有噪声的情况下测试本文算法的性能.其他Actor则负责探索环境,获得更多关于环境的信息,并且与环境交互生成经验数据.Learner每更新 $N$ 步,将网络参数复制给Actor.本文将 $N$ 值设为5.为了加快训练模型以及更好的利用资源,Learner模型在GPU上进行运算,Actor在CPU进行运算.强化学习算法性能评价的两个关键指标是算法收敛所需的步数以及最终收敛的Episode Reward值.收敛所需的步数越少,表示算法的收敛速度越快,最终的收敛的Episode Reward值越大,表示智能体的学习效果越好.本文将根据这两个指标对本文算法进行评价.并且本文还做了一些消融实验,分别删除本文算法中的单项组件,以确定其特定的贡献.

本文实验中统一使用大小为 $1 \times 10^6$ 经验重放池存储智能体与环境交互产生的经验数据.为确保一定的探索性,在智能体的动作策略中添加Ornstein-Uhlenbeck噪声<sup>[19]</sup>.并且本文实验中actor和critic的学习率以及采样批次大小,软更新频率统一设置如下:

$$\alpha = 5 \times 10^{-5}, \beta = 5 \times 10^{-5}, M = 256, TAU = 0.001$$

### 4.1 实验环境配置

本文算法中实验环境的配置如表1所示.本文采用了PyTorch作为深度学习框架,在GPU上使用CUDA加速训练模型.

表1 本文实验环境配置

设备名称	设备信息
CPU	AMD Ryzen 9 5900X, 12c, 24t
GPU	Nvidia GeForce RTX 3090
操作系统	Ubuntu 18.04 LTS
CUDA版本	11.1.0
CUDNN版本	8.1.0
Python版本	3.6.12
PyTorch版本	1.7.1

### 4.2 实验场景

本文实验包含了3个实验场景,如图2所示.这3个实验场景分别是多智能体合作清洁,多无人机监控,多智能体救援.本节将对这3个实验场景进行简单描述.

在多智能体清洁场景中,两个或者多个智能体将合作打扫一个房间,实验场景如图2(a)所示.环境中智

能体的任务是在规定的时间内尽量扩大清洁面积,或者是在最短的时间内尽量清洁完所有的区域.该任务实际上是一个路径规划任务,智能体之间应该合作规划具有最小重叠区域的路径.每个智能体的局部观测由周围8个方块的状态组成,一共8维.动作空间是1维,可以朝上下左右4个方向移动.

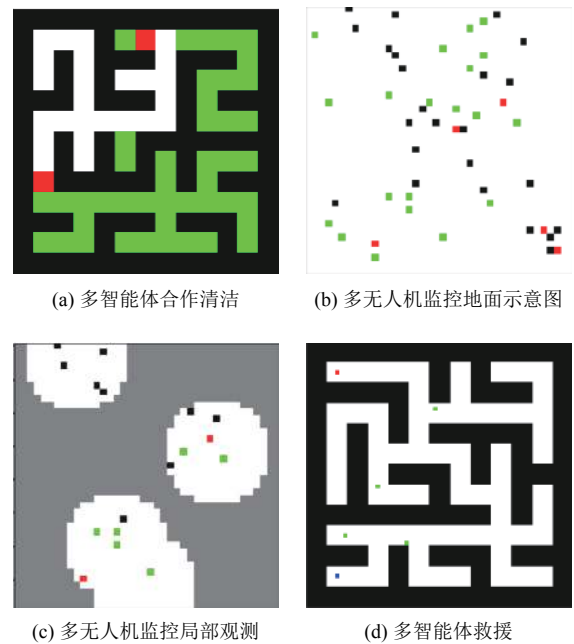


图2 实验场景示意图

多无人机监控实验场景如图2(b)所示.正方形区域表示地面,空中有多个无人机使用摄像头监控地面,无人机可以在给定区域自由飞行.每架无人机都有一个由半径定义的局部观测区域(如图2(c)所示),无人机只能观测到该区域内的物体.在该实验场景中,多个智能体(无人机)的目标是尽量减少视觉区域的重叠以及监控更多的行人.每个智能体的动作空间包含朝上下左右4个方向移动以及停留原地5个动作.

多智能体救援实验场景如图2(d)所示.方形区域模拟地震后的废墟.多个智能体从安全点出发,进入该区域,找到该区域的人并带回安全点.该任务的目标是智能体尽可能在更短的时间内将更多的人带回安全点.每个智能体的局部观测范围是其前方90度的扇形区域,半径为3个方块的长度.智能体的动作空间维度大小为4,分别表示沿 $x$ 轴方向行走 $[-0.5, 0.5]$ ,沿 $y$ 轴方向行走 $[-0.5, 0.5]$ ,旋转 $[-10, 10]$ ,背起/放下人 $\{0, 1, 2\}$ .

3个实验场景中的超参数如表2所示.

表2 本文实验场景的超参数

实验场景	Actor	多步回报- $N$	值分布		
			$N$	$V_{\min}$	$V_{\max}$
多智能体合作清洁	{1,4}	{1,5}	51	-200	400
多无人机监控	{1,4}	{1,5}	51	-600	7000
多智能体救援	{1,4}	{1,5}	51	0	10000

### 4.3 实验结果分析

在框架上, 本文较多借鉴了目前性能最好的多智能体深度强化学习算法-多智能体深度确定性策略梯度算法, 因此选择与该算法进行对比实验. 此外还进行了消融实验, 分别确定本文所提每个方法的贡献. 图3展示了本文算法(MA-D4PG)以及3个消融实验和MA-DDPG在该实验环境中的性能. 其中MAD4PG-no distributional所示曲线表示本文算法去掉第3.1节所述改进后的运行效果, MAD4PG-no multistep所示曲线表示本文算法去掉第3.2节所述改进的运行效果, MAD4PG-no distributed所示曲线表示本文算法去掉第3.3节所述改进的运行效果.

从图3(a)中可看到在多智能体合作清洁的场景中, 本文算法MA-D4PG在收敛速度和最终收敛的Episode Reward值这两个指标上都有最好的表现, 并且算法的稳定性更好. 对MAD4PG-no distributional/MAD4PG-no distributed这两条曲线进行分析可看出, 第3.1节和第3.3节所述的改进对算法最后收敛的Episode Reward值的提升虽不太明显, 但是加快了算法的收敛速度.

对图3(b)进行分析可得, 在多人机监控场景中, 相较MA-DDPG算法, 本文算法(MA-D4PG)收敛速度较慢, 但是最终收敛的Episode Reward值约是MA-DDPG的3.6倍, 智能体的决策能力得到了极大地增强. 对MAD4PG-no distributional/MAD4PG-no distributed分析可得, 第3.1节和第3.3节所述改进使得最终收敛的Episode Reward值分别提升了约97%和43%.

对图3(c)分析可得, 在多智能体救援场景中, 本文算法(MA-D4PG)的收敛速度约是MA-DDPG算法的3.75倍. 对MAD4PG-no multistep/MAD4PG-no distributed两条曲线分析可得, 第3.2节和第3.3节所述改进使得算法的收敛速度分别提升了约185%和100%.

对图3(a), 图3(b)的MAD4PG-no multistep曲线联合分析可看出, 去掉多步回报的改进, 算法性能下降最大, 并且去掉该改进会导致算法的稳定性显著下降.

这两个实验场景中智能体只包含一维的离散动作空间, 控制任务相对简单. 从中可看出在离散控制的多智能体场景中, 引入多步回报对算法性能的贡献最大, 而第3.1节和第3.3节所述改进之间的互补作用较小.

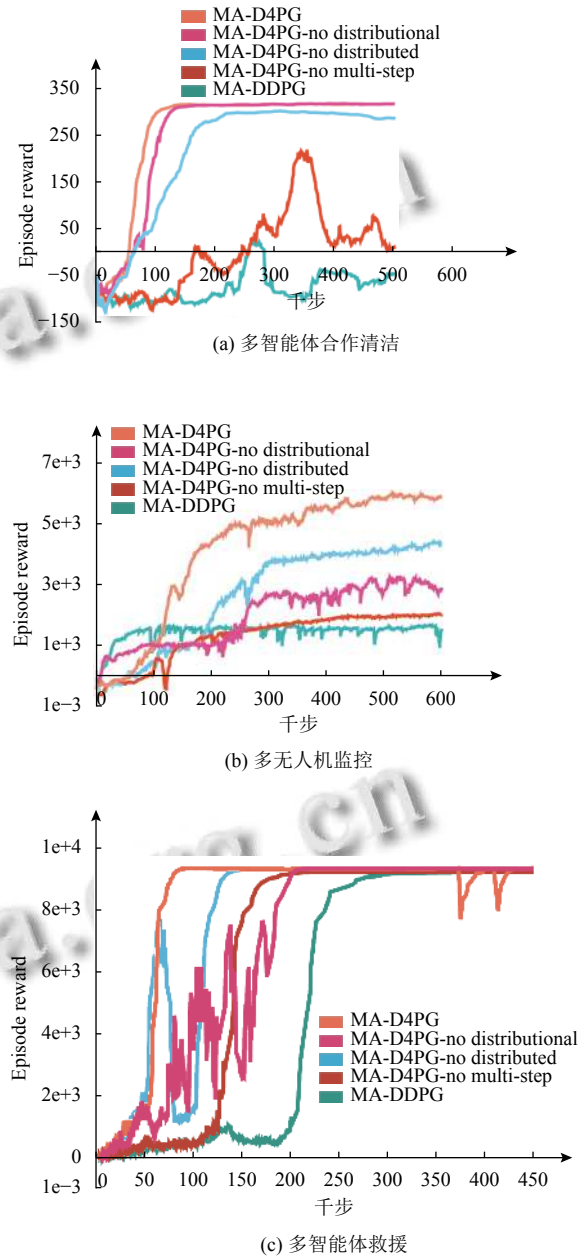


图3 本文算法性能对比图

对图3(c)的MAD4PG-no distributional曲线分析可得, 去掉值分布的改进, 算法的收敛速度降低最大, 并且算法的稳定性也有一定的下降. 该实验环境中的控制任务比较复杂, 既包含连续动作, 又包含离散动作, 而前两个实验场景中智能体只包含一维的离散动作空

间. 从中可看出在复杂的控制任务中, 引入值分布对算法性能的贡献最为突出.

## 5 结论与展望

本文的主要贡献是提出了一种基于值分布的多智能体分布式策略梯度算法, 利用值分布替代期望值, 进一步区分智能体行为的好坏, 而不是简单地将其平均为期望值; 同时引入了多步回报增加智能体的预见能力, 提高算法稳定性; 此外引入了分布式数据生成框架将多智能体场景中智能体与环境交互生成经验数据和网络更新解耦, 加快算法的收敛. 本文算法在3个实验场景中都实现了最好的性能. 这3项改进都有助于提升本文算法(MA-D4PG)的整体性能, 在前两个实验场景中对性能贡献最大的改进是多步回报. 而在复杂的任务场景中, 值分布的贡献变得更加突出. 本文的实验场景都是基于多智能体合作的场景, 接下来的研究工作是在竞争场景甚至既有合作又有竞争的混合场景中研究本文算法的性能.

### 参考文献

- Mnih V, Kavukcuoglu K, Silver D, *et al.* Human-level control through deep reinforcement learning. *Nature*, 2015, 518(7540): 529–533. [doi: [10.1038/nature14236](https://doi.org/10.1038/nature14236)]
- Andrychowicz M, Baker B, Chociej M, *et al.* Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 2020, 39(1): 3–20. [doi: [10.1177/0278364919887447](https://doi.org/10.1177/0278364919887447)]
- Dabney W, Kurth-Nelson Z, Uchida N, *et al.* A distributional code for value in dopamine-based reinforcement learning. *Nature*, 2020, 577(7792): 671–675. [doi: [10.1038/s41586-019-1924-6](https://doi.org/10.1038/s41586-019-1924-6)]
- Kiran BR, Sobh I, Talpaert V, *et al.* Deep reinforcement learning for autonomous driving: A survey. arXiv: 2002.00444, 2021.
- Hua YX, Li RP, Zhao ZF, *et al.* GAN-powered deep distributional reinforcement learning for resource management in network slicing. *IEEE Journal on Selected Areas in Communications*, 2020, 38(2): 334–349. [doi: [10.1109/JSAC.2019.2959185](https://doi.org/10.1109/JSAC.2019.2959185)]
- Liang W, Huang WH, Long J, *et al.* Deep reinforcement learning for resource protection and real-time detection in IoT environment. *IEEE Internet of Things Journal*, 2020, 7(7): 6392–6401. [doi: [10.1109/JIOT.2020.2974281](https://doi.org/10.1109/JIOT.2020.2974281)]
- Chen MJ, Wang T, Ota K, *et al.* Intelligent resource allocation management for vehicles network: An A3C learning approach. *Computer Communications*, 2020, 151: 485–494. [doi: [10.1016/j.comcom.2019.12.054](https://doi.org/10.1016/j.comcom.2019.12.054)]
- Nguyen TT, Nguyen ND, Nahavandi S. Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications. *IEEE Transactions on Cybernetics*, 2020, 50(9): 3826–3839. [doi: [10.1109/TCYB.2020.2977374](https://doi.org/10.1109/TCYB.2020.2977374)]
- Sukhbaatar S, Lin ZM, Kostrikov I, *et al.* Intrinsic motivation and automatic curricula via asymmetric self-play. arXiv: 1703.05407, 2018.
- Palmer G, Tuyls K, Bloembergen D, *et al.* Lenient multi-agent deep reinforcement learning. arXiv: 1707.04402, 2018.
- Foerster J, Farquhar G, Afouras T, *et al.* Counterfactual multi-agent policy gradients. *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*. New Orleans: AAAI, 2018. 2974–2982.
- Lowe R, Wu Y, Tamar A, *et al.* Multi-agent actor-critic for mixed cooperative-competitive environments. arXiv: 1706.02275, 2020.
- Bellemare MG, Dabney W, Munos R. A distributional perspective on reinforcement learning. *Proceedings of the 34th International Conference on Machine Learning*. Sydney: PMLR, 2017. 449–458.
- Horgan D, Quan J, Budden D, *et al.* Distributed prioritized experience replay. arXiv: 1803.00933, 2018.
- Silver D, Lever G, Heess N, *et al.* Deterministic policy gradient algorithms. *Proceedings of the 31st International Conference on Machine Learning*. Beijing: JMLR, 2014. 387–395.
- Lillicrap TP, Hunt JJ, Pritzel A, *et al.* Continuous control with deep reinforcement learning. arXiv: 1509.02971, 2019.
- Barth-Maron G, Hoffman MW, Budden D, *et al.* Distributed distributional deterministic policy gradients. arXiv: 1804.08617, 2018.
- Hessel M, Modayil J, Van Hasselt H, *et al.* Rainbow: Combining improvements in deep reinforcement learning. *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*. New Orleans: AAAI, 2018. 3215–3222.
- Kingma DP, Ba J. Adam: A method for stochastic optimization. arXiv: 1412.6980, 2017.