

基于强化学习的多机协同超视距空战决策算法

王志刚¹, 龚华军¹, 尹逸², 刘小雄²

(1. 南京航空航天大学 自动化学院, 南京 210016)

(2. 西北工业大学 自动化学院, 西安 710072)

摘要: 现代战争中的空战态势复杂多变, 因此探索一种快速有效的决策方法十分重要。本文对多架无人机协同对抗问题展开研究, 提出一种基于 LSTM-MADDPG 的多机协同超视距空战决策算法。首先, 建立无人机运动模型、雷达探测区模型和导弹攻击区模型。然后, 提出了多机协同超视距空战决策算法: 设计了集中式训练分布式执行架构和协同空战系统的状态空间来处理多架无人机之间的同步决策问题; 设计了学习率衰减机制来提升网络的收敛速度和稳定性; 利用长短期记忆网络改进了网络结构, 增强了网络对战术特征的提取能力; 利用基于衰减因子的奖励函数机制加强无人机的协同对抗能力。仿真结果表明所提出的多机协同超视距空战决策算法使无人机具备了协同攻防的能力, 同时算法具备良好的稳定性和收敛性。

关键词: 协同空战决策、多智能体强化学习、混合奖励函数、长短期记忆网络

中图分类号: {V249.1} 文献标识码: A 文章编号: xxxx-xxxx(xxxx)xx-xxxx-xx

Multi aircraft collaborative beyond-visual-range air combat decision-making algorithm based on reinforcement learning

Zhigang Wang¹, Huajun Gong¹, Yi Yin², Xiaoxiong Liu²

(1. Collage of Automation Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China)

(2. School of Automation, Northwestern Polytechnical University, Xi'an 710129, China)

Abstract: With the modern air combat environment is becoming more and more complex, and the combat situation is changing rapidly, so it is important to explore a fast and effective decision-making method. In this paper, we propose an LSTM-MADDPG-based multi-aircraft collaborative beyond-visual-range air warfare decision algorithm to study the problem of multi-aircraft collaborative confrontation. Firstly, a beyond-visual-range air combat environment is established, including UAV movement model, radar detection zone model and missile attack zone model. Then, the multi-aircraft collaborative beyond-visual-range air warfare decision-making algorithm is proposed: a centralized training distributed execution architecture and the state space of the collaborative air warfare system are designed to deal with the synchronous decision-making problem among multiple UAVs; a learning rate decay mechanism is designed to improve the convergence speed and stability of the network; the network structure is improved using a long and short-term memory network to enhance the network's ability to extract tactical features; a decay-based factor-based reward function mechanism to enhance the cooperative countermeasure capability of UAVs. Finally, the results show that the proposed algorithm enables UAVs to have the ability of collaborative attack and defenses, while the algorithm has good stability and convergence.

Keywords: cooperative air warfare decision making, multi-intelligence reinforcement learning, hybrid reward functions, long and short-term memory networks

收稿日期: 2024-xx-xx

作者简介:

通信作者:

1 引言

随着无人机和人工智能技术越来越成熟,无人机具备了无人员伤亡、成本低、持续作战能力强等特点,同时随着人工智能技术的蓬勃发展,未来空战会将进入智能化时代^[1]。多架高性能无人机组成的战术编队体系在空战数据链的加持下相互补充、相互协调进行协同超视距空战。**超视距空战是一种在空中作战中,对抗双方在目视距离(一般为8km)之外探测对方位置,从而使用导弹进行攻击,其特点为作战距离远、多目标攻击、体系对抗、攻防一体;超视距空战与近距空战的主要区别在于作战距离和作战方式。**随着空战环境的复杂多变,空战战场环境愈加错综复杂,飞行员需要处理大量信息。面对多变的战场环境,难以及时感知态势变化,从而对飞机进行飞行决策,并在飞行过程中根据战场环境做出实时调整。所以空战智能决策技术成为了当前研究的重要问题^[2]。

超视距空战智能决策即在飞行员目视距离外,通过人工智能赋能空战对抗,从而辅助或替代飞行员进行空战机动及载荷调度决策的技术^[3],该技术的研究起始于上世纪60年代,该阶段的研究以美国国家航空航天局(National Aeronautics and Space Administration, NASA)兰利研究中心资助的自适应机动逻辑(Adaptive Maneuvering Logic, AML)^[4,5]为代表,其核心为由空战专家总结的空战规则组成的知识驱动型专家系统。到了上世纪90年代,NASA兰利研究中心在AML的基础上,进一步开发了PALADIN系统^[6,7],相比AML,PALADIN系统的最大创新点在于其规则是基于空战仿真自动生成。到了2016年,基于美国空军研究实验室(Air Force Research Laboratory, AFRL)的AFSIM仿真系统开发的“阿尔法空战”系统^[8]战胜了美国退役空军上校基恩·李,与它类似机理的还有波音公司开发的双边对抗学习系统,这两个系统都是首先基于人类经验设计策略结构,随后基于对抗博弈实现参数演进,只是后者的环境适应性更强。2010年之后,基于机器学习的空战决策技术逐渐得到了发展,其中一项广为人知的项目便是美国国防高级研究计划局主导的“阿尔法狗斗”智能近距空中格斗项目,该项目中,苍鹭系统公司的智能决策系统最终以5:0的优势大胜F-16飞行教官Banger^[9]。

由上述叙述可以看出,空战智能决策技术已经发展到了基于强化学习的算法的阶段。通过强化学

习的方法,无人机可以在仿真环境中不断调整自己的空战策略,模拟演练空战中可能遇到的各种情况,形成空战战术。针对空空导弹对空战系统状态多维复杂的问题,张强等人设计了基于Q网络的强化学习和与导弹攻击区相关的奖励函数,形成了一套超视距空战决策系统^[10]。为解决在复杂环境中,飞机进行自主态势决策的问题,李永丰应用改进的深度强化学习算法进行空战任务决策^[11]。朴海音等人则通过多智能体强化学习方法训练出了超越专家水平的超视距空战智能决策系统^[12]。强化学习算法为超视距空战决策问题提出了一种新的解决方案。但是,目前基于强化学习的超视距空战决策方法多研究1V1超视距空战决策,对于多机协同超视距空战决策算法研究较少,且现有协同超视距空战决策算法训练效率较低,花费成本较高^[13-15]。

针对上述问题,本文提出了一种基于多智能体强化学习的多机协同超视距空战算法:针对多架无人机同步决策的问题,设计了集中式训练分布式执行架构和协同空战系统的状态空间;利用学习率衰减机制提升神经网络的收敛性能;针对敌方战术特征提取的问题,采用LSTM网络处理具有时序特征的空战数据,提取空战战术特征;利用基于衰减因子的奖励函数机制来加强无人机的协同对抗能力。最后,通过仿真分析验证了本文所提算法的有效性。

2 超视距空战环境建模

根据多架飞机协同超视距空战决策要求,本文设计了强化学习总体架构,该空战决策结构由智能体模型、无人机运动控制模型、空战环境、空战态势感知、态势评估等组成。

智能体决策部分利用强化学习机制建立了由战场环境信息到无人机控制量的映射,使得无人机可以根据战场实时态势调整飞行轨迹从而完成作战任务;无人机运动控制模型执行智能体决策部分给出的无人机运动控制量,完成无人机的状态更新;敌我双方的无人机、各自的机载火控雷达探测区和空空导弹攻击区共同构成了战场环境;态势感知和评估部分能够感知战场环境的变化,根据环境状态信息利用奖励函数对无人机执行的控制量进行评价。根据强化学习环境的总体架构。

2.1 无人机运动控制模型

为了准确描述无人机的飞行特征,在航迹坐标系下建立无人机的模型:

$$\begin{cases} dv/dt = g(n_x - \sin \theta) \\ d\psi/dt = gn_z \sin \phi / v \cos \theta \\ d\theta/dt = (g/v)(n_z \cos \phi - \cos \theta) \\ dx/dt = v \cos \theta \cos \psi \\ dy/dt = v \cos \theta \sin \psi \\ dz/dt = v \sin \theta \end{cases} \quad (1)$$

式(1)包含六个一阶微分方程定义了无人机的速度 v 、航迹方位角 ψ 、航迹倾斜角 θ 和三轴位置 (x, y, z) 。本文将切向过载 n_x 、法向过载 n_z 和航迹滚转角 ϕ 组合为一个三元组 $[n_x, n_z, \phi]$ 作为智能体决策模块的输出来控制无人机完成一系列的飞行动作。

根据载机和目标的速度与位置信息可以得到载机和目标的相对关系，两者的相对关系可以为无人机的机载火控雷达探测区和导弹攻击区提供判断条件。通过式(2)可以计算得到如图 1 所示的相对关系。

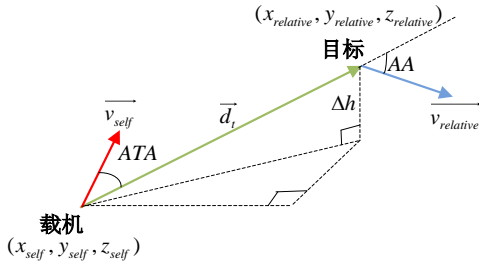


图 1 载机和目标的相对关系图

Figure 1 Relative relationship between the carrier and the target

$$\begin{cases} v_{self} = [\cos \psi_{self} \cos \theta_{self}, \sin \psi_{self} \cos \theta_{self}, \sin \theta_{self}] \\ v_{relative} = [\cos \psi_{relative} \cos \theta_{relative}, \sin \psi_{relative} \cos \theta_{relative}, \sin \theta_{relative}] \\ \vec{d} = [x_{self} - x_{relative}, y_{self} - y_{relative}, z_{self} - z_{relative}] \\ d = |\vec{d}| \\ AA = \arccos(v_{relative} \cdot \vec{d} / d) \\ ATA = \arccos(v_{self} \cdot \vec{d} / d) \\ \beta = \arccos(v_{self} \cdot v_{relative}) \\ \Delta h = z_{self} - z_{relative} \end{cases} \quad (2)$$

其中， v_{self} 为载机的速度方向； $v_{relative}$ 为目标的速度方向； θ_{self} 和 ψ_{self} 为载机的航迹倾斜角和航迹方位角； $\theta_{relative}$ 和 $\psi_{relative}$ 为目标的航迹倾斜角和航迹方位角； x_{self} 、 y_{self} 和 z_{self} 为载机的三轴位置； $x_{relative}$ 、 $y_{relative}$ 和 $z_{relative}$ 为目标的三轴位置； \vec{d} 为载机到目标的距离矢量， d 为载机与目标的距离； AA 为脱离角，即目标的飞行速度与敌我双方距离矢量的夹角； ATA 为偏离角，即载机的飞行速度

与敌我双方距离矢量的夹角； β 为两机速度矢量的夹角； Δh 为载机和目标的高度差。

2.2 机载火控雷达探测区模型

本文根据雷达的属性和特征建立了机载火控雷达的探测区。机载火控雷达在搜索、确认、跟踪 3 个工作阶段的转换关系，每个阶段的主要功能和雷达建模如下：

(1) 空域搜索

在无人机进入空战战场后，进入空域搜索阶段，向指定空域内发射周期性的电磁波进行探测。当电磁波遇到目标时会反射回来，若无人机接收到电磁波回波，则可以确认目标的信息。能否探测到目标与当前我机的航迹姿态、敌我双方的高度和偏离角有关。假设红方为我方，蓝方为敌方，建立如式(3)所示的雷达探测区：

$$\begin{cases} \Delta h_{min} \leq |\Delta h| \leq \Delta h_{max} \\ d_t \leq d_{rmax} \\ -\theta_{rmax} \leq ATA \leq \theta_{rmax} \\ -\psi_{rmax} \leq \psi \leq \psi_{rmax} \end{cases} \quad (3)$$

其中， Δh_{min} 和 Δh_{max} 为 AFR 的最小搜索高度差和最大搜索高度差； Δh 为敌我双方的高度差 $\Delta h = h_t^r - h_t^b$ ， h_t^r 和 h_t^b 为 t 时刻红方和蓝方的高度； d_t 为 t 时刻红方和蓝方的距离； d_{rmax} 为 AFR 的最大探测距离； θ_{rmax} 和 ψ_{rmax} 为 AFR 的最大搜索俯仰角和最大搜索偏航角。

(2) 确认目标

当目标满足式(3)的四个条件时，则符合了被我方发现的先决条件，我方雷达开始确认目标状态。若目标被确认，则其雷达告警系统会报警，表明自身被我方雷达发现，目标越靠近雷达准线，越容易被我方雷达发现；敌我双方的距离越近，越易被我方雷达发现。具体关系如式(4)所示：

$$P_r = (1 - \frac{ATA}{\theta_{rmax}}) \cdot (1 - \frac{ATA}{\psi_{rmax}}) \cdot e^{-\sigma \frac{d_t}{d_{rmax}}} \quad (4)$$

其中， σ 是与 AFR 的散射截面相关的参数，当散射截面积为 $5m^2$ 时， σ 为 0.1625。若 $P_r > 0.2$ ，则成功确认目标，若未能成功确认目标，雷达将重新进入空域搜索阶段，若成功确认目标，雷达将转入跟踪目标阶段。在探测并确认目标后，雷达会进入跟踪状态，此时机载火控系统会计算导弹攻击区和相应的击毁概率，一旦满足打击条件立即发射导弹击毁目标。

(3) 跟踪阶段

在探测并确认目标后，雷达会进入跟踪状态，

此时机载火控系统会计算导弹攻击区和相应的击毁概率，一旦满足打击条件立即发射空空导弹击毁目标。

2.3 导弹攻击区模型

考虑到导弹性能对无人机作战能力的限制，本小节将用衡量中远程空空导弹的战斗能力的六个要素，分别为导弹的最大离轴发射角 φ_{mmax} 、最大攻击距离 d_{mmax} 、最小攻击距离 d_{mmin} 、不可逃逸的圆锥角 φ_{memax} 、不可逃逸的最大距离 d_{memax} 和不可逃逸的最小距离 d_{memin} 来完成攻击区的建模，具体划分如式(5)所示：

$$Area = \begin{cases} Area_{attack}, & d_{mmin} \leq d_t \leq d_{mmax} \text{ and } ATA \leq \varphi_{mmax} \\ Area_{noescape}, & d_{memin} \leq d_t \leq d_{memax} \text{ and } ATA \leq \varphi_{memax} \end{cases} \quad (5)$$

当目标进入导弹攻击区后，根据所处位置不同，被击毁的概率也不同。无人机具备了通过一些连续机动规避导弹攻击的能力，所以在计算导弹击毁概率时，需要考虑目标此时能否通过连续机动规避导弹。本小节通过载机的偏离角 ATA 和目标的脱离角 AA 来描述空战双方的规避优势从而定量分析导弹的击毁概率： ATA 越小，载机的导弹离轴发射角越小，导弹更容易命中目标；而 AA 越小，则目标的飞行方向越接近雷达准线和攻击区中线，越难通过连续机动躲避雷达跟踪和导弹攻击。结合 UAV 的规避优势，将攻击区进一步划分为 5 个部分，每部分的毁伤概率不同。

目标的毁伤概率为：

$$P_{ad} = \tau_a \cdot P_a + \tau_d \cdot P_d \quad (6)$$

其中， P_a 为与目标规避优势相关的毁伤概率， P_d 为与敌我距离相关的毁伤概率， τ_a 和 τ_b 分别为规避优势和距离毁伤概率的权重因子，具体根据双方偏离角、脱离角和距离决定。

$$P_d = \begin{cases} AA/\pi + 1 & \text{if } position(aircraft_aim) \in 1 \\ -(AA/(\pi/6 + ATA)) \times AA + 0.5 & \text{if } position(aircraft_aim) \in 2 \text{ and } \arctan v_r / v_r \geq 0 \text{ and } AA < 5\pi/6 - ATA \\ -(0.5/(5\pi/6 - ATA)) \times AA - 0.5 \times \pi/6 + ATA/(5\pi/6 - ATA) & \text{if } position(aircraft_aim) \in 2 \text{ and } \arctan v_r / v_r \geq 0 \text{ and } AA \geq 5\pi/6 - ATA \\ (0.3/(\pi/6 + ATA)) \times AA + 0.5 & \text{if } position(aircraft_aim) \in 2 \text{ and } \arctan v_r / v_r < 0 \text{ and } AA < 5\pi/6 - ATA \\ (0.3/(ATA - 5\pi/6)) \times AA + (0.5 \times ATA - 43/60)/(ATA - 5\pi/6) & \text{if } position(aircraft_aim) \in 2 \text{ and } \arctan v_r / v_r < 0 \text{ and } AA \geq 5\pi/6 - ATA \\ (0.3/(\pi/6 + ATA)) \times AA + 0.5 & \text{if } position(aircraft_aim) \in 3 \text{ and } \arctan v_r / v_r \geq 0 \text{ and } AA < 5\pi/6 - ATA \\ (0.3/(ATA - 5\pi/6)) \times AA + (0.5 \times ATA - 43/60)/(ATA - 5\pi/6) & \text{if } position(aircraft_aim) \in 3 \text{ and } \arctan v_r / v_r \geq 0 \text{ and } AA \geq 5\pi/6 - ATA \\ -(AA/(\pi/6 + ATA)) \times AA + 0.5 & \text{if } position(aircraft_aim) \in 3 \text{ and } \arctan v_r / v_r < 0 \text{ and } AA < 5\pi/6 - ATA \\ -(0.5/(5\pi/6 - ATA)) \times AA - 0.5 \times (\pi/6 + ATA)/(5\pi/6 - ATA) & \text{if } position(aircraft_aim) \in 3 \text{ and } \arctan v_r / v_r < 0 \text{ and } AA \geq 5\pi/6 - ATA \\ P_d = AA/\pi & \text{if } position(aircraft_aim) \in 4 \end{cases}$$

$$P_d = 1 - ((d_t - d_1) / d_2)^2$$

其中， d_1 和 d_2 为与导弹攻击区相关的距离参数； d_t 为载机和目标的距离。

综上所述，当目标毁伤概率大于阈值时，发射导弹将目标击毁。

3 多机协同超视距空战决策算法设计

本文设计多智能体深度确定策略梯度 (Multi-agent Deep Deterministic Policy Gradient, MADDPG) 算法对无人机协同超视距空战决策进行研究，将集中式训练和分布式执行架构与超视距空战决策算法结合，以 2V2 超视距空战为例，设计了基于 LSTM-MADDPG 多机协同超视距空战决策算法。

3.1 LSTM-MADDPG 算法设计

LSTM-MADDPG 算法由多个 LSTM-DDPG 网络组成，算法的特点如下：

(1) 本文采用集中训练和分布执行的策略。应用集中学习方式训练 Critic 网络，执行动作时，每个个体用独立的 Actor 网络选择动作。Actor 网络用局部观测，Critic 需要全局观测。

(2) 改进经验回放记录数据。为了让网络适用环境，训练中的信息由 $(x, x', a_q, \dots, a_n, r_1, \dots, r_n)$ 组成， $x = (o_1, \dots, o_n)$ 表示每个智能体的观测信息。

(3) 利用所有策略的整体效果对网络进行优化，以提高算法的稳定性和收敛速度。

集中式训练分布式执行架构会对 LSTM-DDPG 算法进行改进，每一个智能体对其他智能体的策略进行函数逼近。将智能体策略网络参数设置为 $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ ，智能体的策略为 $\pi = (\pi_1, \pi_2, \dots, \pi_n)$ ，第 n 个智能体的累积期望奖励如式 (7) 所示。

$$J(\theta_n) = E_{s \sim \rho^{\pi, a_n \sim \pi_{\theta_n}}} \left[\sum_{t=0}^{\infty} \gamma^t r_{n,t} \right] \quad (7)$$

针对 DDPG 算法的确定性策略 μ_{θ_n} ，对累积期望奖励求梯度，梯度公式为式 (8) 所示。

$$\nabla_{\theta_i} J(\mu_i) = E_{x, a \sim D} [\nabla_{\theta_i} \mu_i(a_i | o_i) \nabla_{a_i} Q_i^{\mu}(x, a_1, \dots, a_n) |_{a_i = \mu_i(o_i)}] \quad (8)$$

其中， $x = (\varphi; \dots; o)$ 表示观测向量， $Q_i^{\mu}(x, a_1, \dots, a_n)$ 表示第 i 个智能体的状态-动作函数， D 表示网络训练的经验池，经验池中存放训练数据 $(x, x', a_q, \dots, a_n, r_1, \dots, r_n)$ 。采用集中方式的 Critic 网络其更新计算方式如式 (9) 所示。

$$\begin{cases} L(\theta_i) = E_{x, a, r, x'} [(Q_i^{\mu}(x, a_1, \dots, a_n) - y)^2] \\ y_i = r_i + \gamma \bar{Q}_i^{\mu}(x', a'_1, \dots, a'_n) |_{a'_i = \mu'_i(o_i)} \end{cases} \quad (9)$$

其中， \bar{Q}_i^{μ} 表示目标网络， $\mu' = [\mu'_1, \dots, \mu'_n]$ 为目标策略具有滞后更新的参数 θ'_j 。

在算法设计中 Critic 网络采用全局信息进行网络学习，Actor 网络采用局部观测进行网络学习。代价函数为式 (10) 所示。

$$\begin{cases} L(\phi_i^j) = -E_{o_j, a_j} [\log \hat{\mu}_{\phi_i^j}(a_j | o_j) + \lambda H(\hat{\mu}_{\phi_i^j})] \\ L(\theta_i) = E_{x, a, r, x'} [(Q_i^\mu(x, a_1, \dots, a_n) - y)^2] \\ y_i = r_i + \gamma \bar{Q}_i^\mu(x', a_1', \dots, a_n') |_{a_j' = \mu_j^i(o_j)} \end{cases} \quad (10)$$

只要最小化代价函数，就能得到其他智能体策略的逼近，因此上式的 y 可以替换为式 (11)。

$$\begin{cases} L(\theta_i) = E_{x, a, r, x'} [(Q_i^\mu(x, a_1, \dots, a_n) - y)^2] \\ y_i = r_i + \gamma \bar{Q}_i^\mu(x', \hat{\mu}_{\phi_i^1}^1(o_1), \dots, \hat{\mu}_{\phi_i^n}^n(o_n)) \end{cases} \quad (11)$$

针对强策略很难去适应新的对手策略问题，LSTM-MADDPG 应用了一种策略集合的思想，即第 i 个智能体的策略 μ_i 由一个具有 k 个子策略的集合构成，在每一个训练 episode 中只用一个子策略（简称为 $\mu_i^{(k)}$ ）。对每一个智能体，最大化其策略集合的整体奖励为式 (12) 所示。

$$J_e(\mu_i) = E_{k \sim \text{unif}(1, K), s \sim \rho^\mu, a \sim \mu_i^{(k)}} [\sum_{t=0}^{\infty} \gamma^t r_{i,t}] \quad (12)$$

可以为每一个子策略构建一个记忆存储单元 $D_i^{(k)}$ ，去优化策略集合的整体效果，因此对每一个子策略的更新梯度为式 (13) 所示。

$$\nabla_{\theta_i^{(k)}} J_e(\mu_i) = \frac{1}{K} E_{[x, a]_{k \sim D_i^{(k)}}} [\nabla_{\theta_i^{(k)}} \mu_i^{(k)}(a_i | o_i) \nabla_{a_i} Q_i^\mu(x, a_1, \dots, a_n) |_{a_i = \mu_i^{(k)}(o_i)}]$$

表 1 全局状态量设计

Table 1 Global state volume design

状态分量	变量	描述	维数
本机	<i>state_self</i>	本机信息，包括状态、速度、高度、航迹俯仰角和航迹偏航角	5
	<i>state_relative</i>	与友机和敌机的相对信息，包括存活状态、脱离角、偏离角、相对距离和两机速度矢量的夹角	5×3
	<i>state_radar</i>	雷达对两架敌机的状态和雷达告警器对两架敌机的状态	4
	<i>state_missile</i>	导弹对两架敌机的状态和导弹告警器对两架敌机的状态	4
友机	/	与本机类似，但是不含 <i>state_relative</i> 信息	13
敌 1 号机	/	与本机类似，但是不含 <i>state_relative</i> 信息	13
敌 2 号机	/	与本机类似，但是不含 <i>state_relative</i> 信息	13

在集中式训练分布式执行框架中，每一个智能体无人机都有自己的价值网络，价值网络输入即全局状态量 s 均按本机、友机、敌 1 号机和敌 2 号机的排序。私有观测状态量 o_i 也是由本机、友机、敌 1 号机和敌 2 号机的顺序组成，因为友机之间可以通过数据链进行通信，所以私有观测状态量 o_i 中可以包含全部的友机信息，在实际空战中，敌机的雷达和导弹攻击区状态均可以通过本机的雷达告警器和导弹告警器得到，所以在私有观测状态量 o_i 中敌 1 号机和敌 2 号机仅保留自身的状态信息。最后，私有观测状态量 o_i 总共 51 维。

3.3 学习率衰减机制

对于学习率的设置，本文希望在训练前期通过较大的学习率加快网络收敛速度，在训练后期则通

$$\begin{aligned} & \text{计算 TD 误差:} \\ \delta_i^j &= Q(S_i, [\mu(o_i^1 | \theta_i^\mu), \mu(o_i^2 | \theta_i^\mu), \mu(o_i^3 | \theta_i^\mu), \mu(o_i^4 | \theta_i^\mu)] | w_i^Q) - \hat{y}_i^j \quad (14) \end{aligned}$$

然后通过梯度下降算法更新参数 w_i^Q 使得价值网络的预测值更接近 TD 目标值：

$$w_i^Q = w_i^Q - \alpha \cdot \delta_i^j \cdot \nabla_{w_i^Q} Q(S_i, [\mu(o_i^1 | \theta_i^\mu), \mu(o_i^2 | \theta_i^\mu), \mu(o_i^3 | \theta_i^\mu), \mu(o_i^4 | \theta_i^\mu)] | w_i^Q) \quad (15)$$

最后可以通过软更新的方法更新每一架无人机的目标策略网络和目标价值网络：

$$\begin{cases} \theta_i^{\mu'} = \tau \theta_i^\mu + (1 - \tau) \theta_i^{\mu''} \\ w_i^{Q'} = \tau w_i^Q + (1 - \tau) w_i^{Q''} \end{cases} \quad (16)$$

3.2 协同空战系统的状态空间设计

由于协同空战中智能体大幅增加，根据多机协同超视距空战和集中式训练分布式执行架构的特点，设计了协同空战系统的全局状态量 s 和私有观测状态量 o_i 。以 2V2 协同超视距空战为例，智能体有 4 个，分别为本机、友机、敌 1 号机和敌 2 号机，其全局状态量 s 如表 1 所示。

过较小的学习率来找到全局最优点。为实现这一目标，本小节设计了如式(17)的学习率衰减机制。

$$l_{decayed} = l_0 \cdot (m^{\frac{N_{episode}}{N_{decay}}}) \quad (17)$$

其中， $l_{decayed}$ 为衰减后的学习率； l_0 为初始设置的学习率； m 为衰减比例； $N_{episode}$ 为总的训练回合； N_{decay} 为衰减步数。

在超视距空战决策算法训练过程中，根据本节设计的学习率衰减机制，学习率可以随着训练回合而变化。在算法训练初期，学习率较大，可以让无人机的空战策略迅速向最优策略靠拢，随着训练的进行，学习率逐渐减小，降低网络梯度变化的速度，使得无人机学习的空战策略不会错过最优策略。学习率衰减机制既保证了超视距空战决策算法的训练速度，又避免了算法陷入局部最优。

3.4 混合奖励函数及其校正机制

本小节设计的混合奖励包括状态奖励和事件奖励，并且根据多机协同超视距空战的特点，增加了回合奖励和衰减因子。

(1) 状态奖励

状态奖励的目的是使无人机能够学会安全飞行并且引导无人机向目标靠近，产生与目标对抗的经验，所以本小节通过引导奖励来实现这一目的。因为在多机协同超视距空战中存在多个目标，所以设计了一个目标分配方法，先计算我方两架无人机与敌方两架无人机之间的相对距离，然后选择相对距离最短的两架无人机相互作为目标，而剩下的两架无人机则互为目标，所以引导奖励 r_{state} 如下：

$$r_{state} = \begin{cases} 0.01 & \Delta d > 50 \\ 0 & -50 \leq \Delta d \leq 50 \\ -0.01 & \Delta d < -50 \end{cases} \quad (18)$$

其中， Δd 为一个决策周期内本机与对应目标之间相对距离的变化量，可由式(19)计算可得。

$$\Delta d = d_{last} - d_{current} \quad (19)$$

其中， d_{last} 为上一个决策时刻本机与对应目标之间相对距离； $d_{current}$ 为当前决策时刻本机与对应目标之间相对距离。

(2) 事件奖励

事件奖励通过给超视距空战中涉及的标志性事件单独设计奖励，从而引导无人机在攻击目标的同时，能够规避敌方的导弹攻击区。主要事件和对应的奖励如表 2 所示。

表 2 事件奖励

Table 2 Event Rewards

事件类型	奖励函数	事件类型	奖励函数
目标进入己方 AFR 搜索区	$r_{event_ afr1}$	己方进入目标 AFR 搜索区	$-r_{event_ afr1}$
己方 AFR 确认目标	$r_{event_ afr2}$	己方被目标 AFR 确认	$-r_{event_ afr2}$
己方 AFR 跟踪目标	$r_{event_ afr3}$	己方被目标 AFR 跟踪	$-r_{event_ afr3}$
目标进入己方导弹攻击区	$r_{event_ missile1}$	己方进入目标导弹攻击区	$-r_{event_ missile1}$
己方导弹未命中目标	$r_{event_ missile2}$	敌方导弹未命中己方	$-r_{event_ missile2}$
己方导弹命中目标	$r_{event_ missile3}$	敌方导弹命中己方	$-r_{event_ missile3}$
超出战场边界	$r_{event_ battlefield}$	飞机速度限制	$r_{event_ v}$

(3) 回合奖励

多机协同超视距空战中，需要将敌方所有无人

机全部击毁才能结束本回合的空战。所以本节设计了回合奖励 $r_{episode}$ ，即在空战回合结束时，获胜的一方智能体可以获得奖励，即 $r_{episode} = 10$ ，而战败的一方智能体则给予惩罚，即 $r_{episode} = -10$ ，而在规定的时间内，空战双方未能分出胜负，双方智能体既不获得奖励也不给予惩罚 $r_{episode} = 0$ 。综上所述，混合奖励函数计算方式如式(20)所示。

$$r = r_{episode} + \max(0, (1 - episode / \lambda_{state})) \cdot r_{state} + \max(0, (1 - episode / \lambda_{event})) \cdot r_{event} \quad (20)$$

其中， λ_{state} 为状态奖励衰减因子， λ_{event} 为事件奖励的衰减因子。

衰减因子的引入，可以在多机协同超视距空战算法前期引导智能体无人机快速学会安全飞行并且积攒与敌方交战的经验，加速算法的收敛。随着训练的增加，获得奖励会只剩下回合奖励，防止算法陷入局部最优，鼓励智能体无人机尝试更多的战术配合，使多机协同超视距空战涌现出更多的战术。

3.5 多机协同超视距空战决策算法框架

基于多机超视距空战的任务场景，本文提出基于 LSTM-MADDPG 的多智能体空战决策框架，该框架主要由三部分组成：深度强化学习模块、环境模块和数据处理模块。具体结构如图 2 所示（以 2V2 情况为例）。强化学习模块主要分为两部分，中央控制器和分布式执行，中央控制器主要负责评估智能体无人机执行的动作；分布式执行主要输出红方无人机和蓝方无人机需要执行的动作。环境模块主要由无人机模型、导弹攻击区和奖励模块组成；奖励模块根据战场态势进行态势评估；空战数据处理模块的功能是处理环境模块发送的空战数据，对空战数据进行掩码处理、归一化和打包，并将其发送到共享经验池。

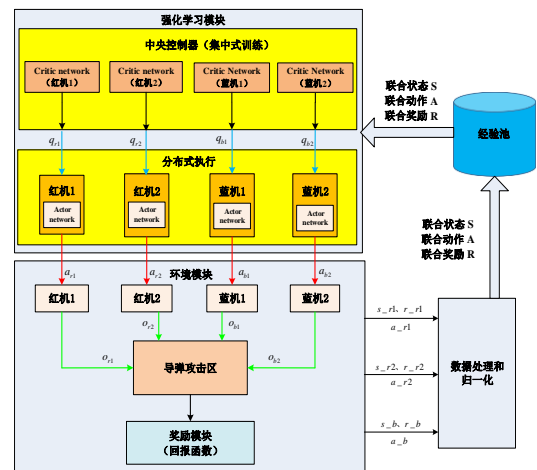


图 2 多机协同超视距空战决策框架图

Figure 2 Multi-aircraft cooperative over-the-horizon air combat decision-making framework diagram

红方无人机和蓝方无人机根据其策略网络的初始参数选择初始动作，红蓝双方的无人机执行此动作与环境交互，以获得新的状态和奖励。此时，数据处理模块对双方的状态、动作和奖励等数据进行打包、归一化、掩码处理，将其处理成联合动作集合、联合奖励集合和联合状态集合。待共享经验池满后，强化学习模块开始对数据进行采样，并将采样的联合状态、联合动作和联合奖励发送给中央控制器，对价值网络进行更新。在价值网络更新完成之后，会对红方和蓝方执行的动作进行评价，指导策略网络完成更新。策略网络更新后，红方和蓝方飞机将自身的观测结果输入网络，策略网络输出相应的动作。无人机在收到强化学习模块的动作决策后，执行相应的动作，与空战环境进行交互，产生新一轮的状态、动作和奖励。数据处理模块将新的数据处理完成后，将其发送到共享经验池中。如此

循环往复，直至满足多机协同超视距空战训练结束的要求。

4 实验结果与分析

为验证本章提出的多机协同超视距空战决策算法的有效性，在 2V2 超视距空战的背景下，进行 LSTM-MADDPG 与 QMIX 和 VDN 两种多智能体强化学习算法的对比实验，验证所设计的优化机制对多机协同超视距空战决策算法训练效果的提升作用；最后，设计 2V2 的超视距空战实验场景，对训练过程中涌现的战术进行分析。

4.1 仿真环境设置

多机协同超视距空战决策算法的神经网络参数如表 3 所示，算法训练参数如表 4 所示。

表 3 神经网络参数设置

Table 3 Neural network parameter settings

神经网络名称	感知层细胞数量	拟合层神经网络结构	激活函数	输出层激活函数
价值网络	512	(512, 256, 64, 1)	Relu	/
策略网络	512	(512, 256, 64, 3)	Relu	tanh

表 4 训练参数设置

Table 4 Training parameter settings

训练参数	数值	训练参数	数值
策略网络的学习率	0.01	价值网络的学习率	0.01
折扣率	0.95	批大小	128
软更新权重	0.01	探索初始能力因子	1
探索能力衰减因子	0.00002	学习率衰减比例	e
学习率衰减步数	2000	状态奖励的衰减因子	3,000
事件奖励的衰减因子	6,000		

4.2 对比实验

本小节设计仿真 2V2 超视距空战，进行 LSTM-MADDPG 与 QMIX 和 VDN 算法的对比实验。在 2V2 超视距空战中有 4 架智能体无人机，由于每架无人机的平均奖励收敛情况基本一致，所以选取其中一架无人机的平均奖励来展示算法的收敛情况，下面将对 3 种算法的训练情况进行分析，如图 3 所示。

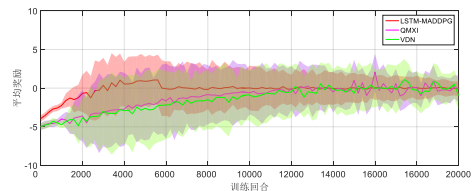


图 3 多智能体强化学习算法平均奖励对比图

Figure 3 Comparison of average rewards of multi-intelligent body reinforcement learning algorithms

首先分析 LSTM-MADDPG 平均奖励的变化情况。在算法训练初期，作为智能体的无人机尚未学会安全飞行，会出现坠毁和飞出空战区域的情况，此时 200 个回合奖励的标准差还不是很大。随着训

练回合的增加，无人机学会了安全飞行，所获得的平均奖励逐渐增加。而双方又在状态奖励的作用下开始接近，随着双方距离的接近，无人机触发了空战中的事件奖励，所以 1600 回合后 200 个回合奖励的标准差开始逐渐变大。6000 回合后，为了防止无人机陷入专家经验，状态奖励和事件奖励衰减为 0，无人机获得的奖励只有回合奖励。随着训练回合的继续增加，无人机开始学会相互协同，进攻脱离或者引诱迂回等战术，平均奖励开始趋于稳定，200 个回合奖励的标准差逐渐减少，LSTM-MADDPG 算法逐渐收敛。

VDN 算法训练收敛趋势与 LSTM-MADDPG 算法类似，但是 VDN 算法的收敛速度远不如 LSTM-MADDPG 算法。QMIX 算法训练收敛趋势也与 LSTM-MADDPG 算法类似，其收敛速度和效果介于 LSTM-MADDPG 算法和 VDN 算法之间。

算法训练完成后，以 2V2 超视距空战为任务场景，调用 VDN 和 QMIX 分别与 LSTM-MADDPG 进行 100 次对抗，对抗结果如图 4 所示。

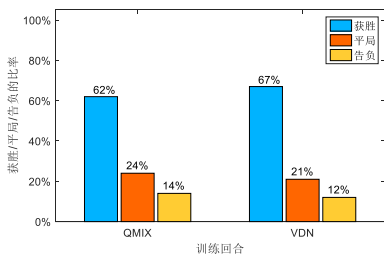


图 4 采用不同算法的空战对抗胜率统计图

Figure 4 Statistical graph of air combat confrontation win rate using different algorithms

若对抗双方同时使用 LSTM-MADDPG 算法，训练过程中胜率变化如图 5 所示。

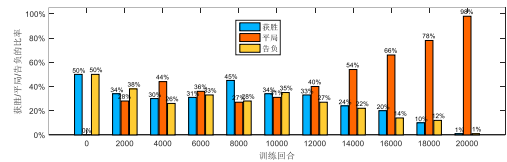


图 5 采用相同算法的空战对抗胜率统计图

Figure 5 Statistical graph of air combat confrontation win rate using the same algorithm

在训练初期，红方和蓝方均没有学会安全飞行，所以此时的告负基本是超出作战区域且红方的获胜和告负概率均为 50%。随着训练的增加，红方和蓝方学会安全飞行，开始尝试对抗，所以红方有胜有负，此时告负中有被蓝方击落的情况。随着持续训练，红方和蓝方激烈对抗，平局的占比减少，获胜和告负的情况增加。在 12000 回合后，红方和蓝方开始学习如何进攻、逃脱对方雷达锁定和导弹攻击区的策略，所以平局的情况开始增加，获胜和告负的情况减少。最后，红方和蓝方在规定的空战时间内，哪一方都无法消灭对方，所以平局的占比达到 98%，通过仿真表明本文设计的空战决策算法满足要求。

4.3 2V2 空战涌现的战术分析

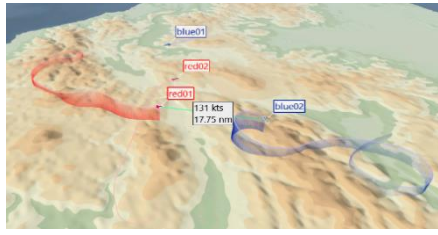
本节设计 2V2 超视距空战实验场景，展示多机协同超视距空战决策算法涌现的空战战术。本小节根据红蓝双方的对抗轨迹、雷达状态和导弹攻击区状态分析对抗双方决策行为的合理性，验证了多架飞机进行协同空战决策的有效性，具体空战场景和红蓝方初始状态设置如表 5 所示。

表 5 红蓝方初始状态设置

Table 5 Red and Blue Initial State Settings

场景	阵营	无人机状态
场景 1: 红蓝方均势 双方迎头飞行	红方 1 号机	(100m / s, 0deg, 0deg, -55km, 50km, 3km)
	红方 2 号机	(100m / s, 0deg, 0deg, -55km, -50km, 3km)
	蓝方 1 号机	(100m / s, 0deg, 180deg, 55km, 50km, 3km)
	蓝方 2 号机	(100m / s, 0deg, 180deg, 55km, -50km, 3km)
场景 2: 红蓝方双方迎头飞行 (不同高度)	红方 1 号机	(100m / s, 0deg, 0deg, -50km, -75km, 2km)
	红方 2 号机	(100m / s, 0deg, 0deg, -50km, -50km, 2km)
	蓝方 1 号机	(100m / s, 0deg, 180deg, 50km, 50km, 3km)
	蓝方 2 号机	(100m / s, 0deg, 180deg, 50km, 75km, 3km)

场景 1 和场景 2 的红蓝双方对抗三维轨迹如图 6 所示。



(a) Scenario 1 3D Air Combat Tracks

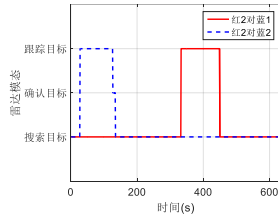


(b) Scenario 2 3D Air Combat Tracks

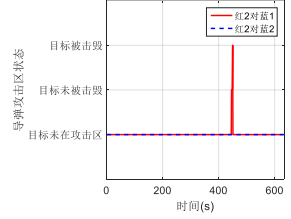
图 6 进攻脱离战术

Figure 6 Offense disengagement tactics

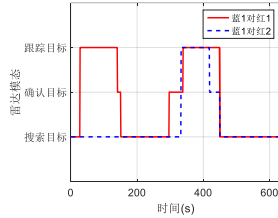
下面以场景 1 为例将涌现出的进攻脱离战术分解为 3 个时刻具体分析，空战轨迹、红蓝双方的雷达和导弹攻击区状态具体如图 7 所示。



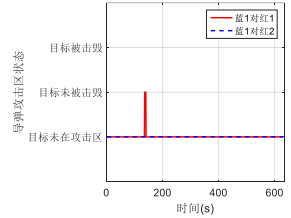
(g) Red 2 radar modes



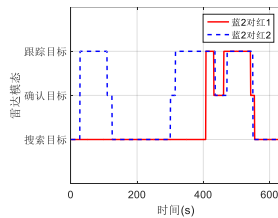
(h) Red 2 missile strike zone



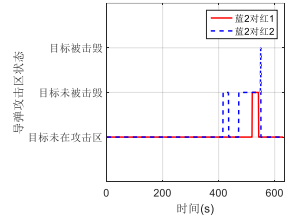
(i) Blue 1 radar modes



(j) Blue 1 missile strike zone



(k) Blue 2 radar modes



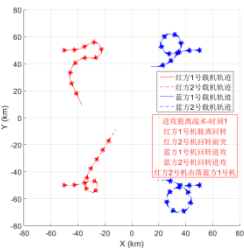
(l) Blue 2 missile strike zone

图 7 进攻脱离战术

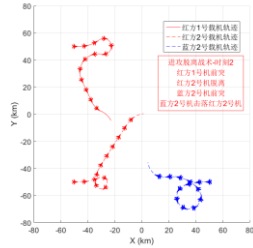
Figure 7 Offense disengagement tactics

如图 7 (a) 所示，红方无人机和蓝方无人机根据目标分配原则，红方 1 号与蓝方 1 号互为目标，红方 2 号和蓝方 2 号互为目标。红方 1 号和蓝方 1 号相互发现后，双方的雷达均转为跟踪模式。红方 1 号与蓝方 1 号受到对方威胁选择了脱离机动，红方 1 号选择掉头躲避威胁，向红方 2 号方向靠拢，而蓝方 1 号在空中完成大转弯后，立刻回转，对红方 1 号发动二次攻击。红方 2 号和蓝方 2 号发现对方后，执行脱离战术。红方 2 号在完成躲避蓝方 2 号攻击后，立刻向蓝方 1 号快速突进，迅速跟踪目标。蓝方 1 号正在向红方 1 号发起第 2 次进攻，当蓝方 1 号受到警告来不及做出脱离机动就被红方 2 号纳入导弹攻击区内，经过红方 2 号的导弹攻击区解算，判定蓝方 1 号被击毁。与此同时蓝方 2 号机成功脱离红方 2 号机的跟踪，立刻回转，向红方 2 号机发起第 2 次进攻，由于自身规避威胁时转弯半径过大，无法有效牵制红方 2 号机。

如图 7 (b) 所示，蓝方 2 号机抓住红方 2 号机攻击蓝方 1 号机的时机，迅速调整攻击角度，进行攻击占位，迅速发现红方 2 号机，转入跟踪模式。而红方 2 号机立刻做出脱离机动，但是由于进攻蓝方 1 号机时太过前突，无法有效躲过蓝方的攻击，

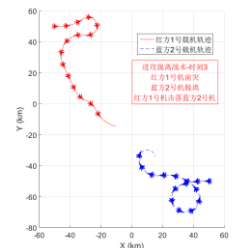


(a) Combat trajectory at 350s

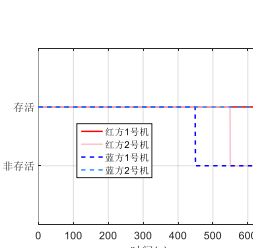


(b) Combat trajectory at 450s

450s

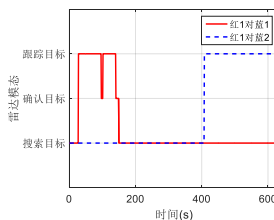


(c) Combat trajectory at 530s

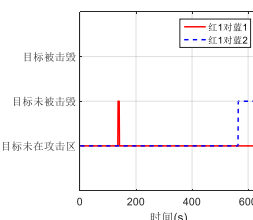


(d) Survival status of parties

parties



(e) Red 1 radar modes



(f) Red 1 missile strike zone

根据蓝方 2 号机的导弹攻击区解算, 判定红方 2 号机被蓝方 2 号机击落。与此同时, 红方 1 号机发现友机受到威胁后, 加速向前前突, 调整好自身的攻击角度进行攻击占位, 对蓝方 2 号机发起进攻。

如图 7 (c) 所示, 红方 1 号机在蓝方 2 号机进攻红方 2 号机时, 对蓝方 2 号机进行攻击占位, 随着两机的距离越来越近, 红方 1 号机的雷达已经跟踪上蓝方 2 号机, 蓝方 2 号机接收到雷达和导弹告警后, 立刻做出脱离机动, 掉头向后跑去。而红方调整好攻击角度后, 加速向前将蓝机纳入导弹攻击区后, 红方 1 号机根据双方态势和蓝机规避优势导弹攻击区解算杀伤概率, 最后判定蓝方 2 号机被红方 1 号机击落。经过一系列对抗, 最终红方取得了本轮 2V2 超视距空战的胜利。

5 总结

本文针对多架飞机, 研究了一种超视距协同空战决策算法。

(1) 根据超视距空战的特点, 给出无人机运动控制模型; 确定机载火控雷达探测区模型及其工作阶段; 根据目标的规避优势、敌我双方距离和双方高度差建立了导弹攻击区和毁伤概率模型。

(2) 采用集中式训练分布式执行架构处理多架无人机同步决策和无人机之间既有协作又有竞争的问题; 针对如何设置学习率的问题, 设计了学习率衰减机制来提升网络的收敛速度和稳定性; 然后, 利用长短期记忆网络改进了网络结构, 使网络可以通过敌我历史动作序列提取敌方的战术特征, 从而做出更有利的空战决策; 提出了基于衰减因子的奖励函数机制和加速网络训练。

(3) 仿真结果表明所提出的多机协同超视距空战决策算法满足无人机协同作战的需求, 仿真中涌现出了一些专家级的超视距协同空战战术。

参考文献:

[1] 杨伟. 关于未来战斗机发展的若干讨论[J]. 航空学报, 2020, 41(6): 524377.
[2] Stillion J. Trends in air-to-air combat implications for

future air superiority[M]. Washington, D.C.: Center for Strategic and Budgetary Assessments, 2015.
[3] 孙智孝, 杨晟琦, 朴海音, 等. 未来智能空战发展综述[J]. 航空学报, 2021, 42(8): 525799.
[4] Burgin G H, Owens A J. An adaptive maneuvering logic computer program for the simulation of one-on-one air-to-air combat[R]. NASA-CR-2582, CR-2583, 1975.
[5] Burgin G. Improvements to the adaptive maneuvering logic program[R]. NASA-CR-3985, 1986.
[6] Goodrich K, Mcmanus J. Development of a tactical guidance research and evaluation system (TGRES)[C]. Boston, MA: Flight Simulation Technologies Conference and Exhibit, AIAA, 1989:
[7] Mcmanus J, Goodrich K. Application of Artificial Intelligence (AI) programming technique stotactical guidance for fighter aircraft[C]. Boston, MA: Guidance, Navigation and Control Conference, AIAA, 1989.
[8] Ernest N, Carroll D. Genetic fuzzy based artificial intelligence for unmanned combat aerial vehicle control in simulated air combat missions[J]. Journal of Defense, 2016, 06(1). DOI: 10.4172/2167-0374.1000144.
[9] Defense Advanced Research Projects Agency. AlphaDogfight trials go virtual for final event[EB/OL]. <https://www.darpa.mil/news-events/2020-08-07>. 2020[2022-04-31].
[10] 张强, 杨任农, 俞立新, 等. 基于 Q-network 强化学习的超视距空战机动决策[J]. 空军工程大学学报(自然科学版). 2018, 19(6): 8-14.
[11] Li Yongfeng, Shi Jingping, Jiang Wei, et al. Autonomous maneuver decision-making for a UCAV in short-range aerial combat based on an MS-DDQN algorithm[J]. Defence Technology. 2022, 18(9): 1697-1714.
[12] Sun Zhixiao, Piao Haiyin, Yang Zhen, et al. Multi-agent hierarchical policy gradient for Air Combat Tactics emergence via self-play[J]. Engineering Applications of Artificial Intelligence, 2021, 98: 104-112.
[13] Isci H, Koyuncu E. Reinforcement Learning Based Autonomous Air Combat with Energy Budgets[C]. USA: AIAA SciTech Forum: AIAA SciTech Forum, 2022: DOI:10.2514/6.2022-0786.
[14] Pan Qian, Zhou Deyun, Huang Jichuan, et al. Maneuver decision for cooperative close-range air combat based on state predicted influence diagram[C]. IEEE, 2017: .DOI:10.1109/ICInfA.2017.8079001.
[15] Wang Luhe, Hu Jinwen, Xu Zhao, et al. Autonomous maneuver strategy of swarm air combat based on DDPG[J]. Journal of Artificial Intelligence and Technology, 2021, 1(1) : 232-243.