

文章编号: 1674-8190(2023)02-145-07

基于近端策略优化的空战决策算法研究

张博超, 温晓玲, 刘璐, 张雅茜, 王宏光

(中国航空工业集团有限公司 沈阳飞机设计研究所, 沈阳 110035)

摘要: 面对未来有/无人机协同作战场景, 实时准确的空战决策是制胜的关键。复杂的空中环境、瞬变的态势数据以及多重繁琐的作战任务, 使有/无人机协同作战将替代单机作战成为未来空战的发展趋势, 但多智能体建模和训练过程却面临奖励分配困难、网络难收敛的问题。针对 5v5 有/无人机协同的空战场景, 抽象出有人机和无人机智能体的特征模型, 提出基于近端策略优化算法的空战智能决策算法, 通过设置态势评估奖励引导空战过程中有/无人机智能体的决策行为向有利态势发展, 实现在与环境的实时交互中, 输出空战决策序列。通过仿真实验对所提空战决策算法进行验证, 结果表明: 本文提出的算法在经过训练学习后, 能够适应复杂的战场态势, 在连续动作空间中得到稳定合理的决策策略。

关键词: 空战决策; 智能决策; 强化学习; 近端策略优化; 有/无人机协同

中图分类号: V212.1; E91

文献标识码: A

DOI: 10.16615/j.cnki.1674-8190.2023.02.17

Research on air combat decision algorithm based on proximal policy optimization

ZHANG Bochao, WEN Xiaoling, LIU Lu, ZHANG Yaqian, WANG Hongguang

(Shenyang Aircraft Design and Research Institute, Aviation Industry Corporation of China, Ltd., Shenyang 110035, China)

Abstract: Facing the future combat scenario with manned and unmanned aerial vehicle cooperation, real-time and accurate air combat decision-making is the basis of winning. The complex air environment, transient situation data, and multiple cumbersome combat tasks make coordinated combat with unmanned aerial vehicles a trend in future air combat, replacing single machine combat. However, multi-agent modeling and training processes face difficulties in reward allocation and network convergence. Air combat scenarios for 5v5 manned and unmanned aerial vehicle cooperation, the characteristic model of single agent is abstracted in this paper, and an algorithm based on proximal policy optimization is proposed to obtain the air combat decision sequence by using reward and punishment incentive in the real-time interaction with the environment. The simulation results show that the algorithm proposed in this paper can adapt to the complex battlefield situation and get a stable and reasonable decision-making strategy in continuous action space after training and learning.

Key words: air combat decision; intelligent decision; reinforcement learning; proximal policy optimization; manned and unmanned aerial vehicle cooperation

收稿日期: 2022-06-11; 修回日期: 2022-09-18

通信作者: 刘璐, aileenliulu@foxmail.com

引用格式: 张博超, 温晓玲, 刘璐, 等. 基于近端策略优化的空战决策算法研究[J]. 航空工程进展, 2023, 14(2): 145-151.

ZHANG Bochao, WEN Xiaoling, LIU Lu, et al. Research on air combat decision algorithm based on proximal policy optimization [J]. Advances in Aeronautical Science and Engineering, 2023, 14(2): 145-151. (in Chinese)

0 引言

面对复杂的空中环境、瞬变的态势数据以及多重繁琐的作战任务,有/无人机协同作战将替代单机作战成为未来空战的发展趋势。协同作战过程中,实时准确的空战决策是制胜的基础。空战决策研究方法可以被分为基于规则和基于智能算法两类。基于规则的决策策略通过态势状态提取指定特征,评估敌我双方的攻击威胁能力,再根据既定规则进行决策,主要包括专家系统^[1]、微分对策算法^[2-3]和矩阵博弈算法^[4-5]等,具有规则固定、人工推演繁琐、无法覆盖瞬时变化的所有状态的特点。基于智能算法的决策策略对战场态势和飞机动作行为进行建模,通过行为与环境的交互结果优化决策模型的结构和参数,主要包括深度学习算法^[6-7]、启发式学习算法^[8-9]、遗传算法^[10]和迁移学习算法^[11]等,具有复杂态势应对性好、自适应强、无需编写明确规则控制的特点,已逐渐成为研究热点。左家亮等^[8]在 2v4 空战场景下,采用启发式强化学习算法提升决策序列的搜索效率;吴宜珈等^[12]提出基于 Opiton 的近端策略分层优化算法,提升空战决策效率,并在 2v2 场景下进行了仿真验证。

针对有/无人机协同、多机协同作战模式下的更多智能体建模和训练过程仍面临奖励分配困难、网络难收敛的问题。本文以一架有人机与四架无人机的 5v5 协同作战为背景,分别建立 5 个智能体,采用强化学习算法^[13]中的近端策略优化算法(Proximal Policy Optimization,简称 PPO),设置智能体间的交互式奖励策略,分别训练决策模型,以获取最优决策。

1 基于近端策略优化的空战决策模型

1.1 近端策略优化算法

强化学习算法包括基于价值函数的学习方法^[14]、基于策略梯度的学习方法^[15]和基于“演员—评论家”(Actor-Critic,简称 A-C)算法^[16]。A-C 算法相当于在基于策略梯度的 Actor 上增加一个基于价值函数进行策略评估的 Critic,此种方式能够增加高价值动作的选择概率,提高智能体的学习速度。由于 5v5 有/无人机协同空战场景复杂,

考虑到实际工程应用中强化学习算法普遍存在的难训练、难收敛等问题,本文选择采用基于 A-C 算法的 PPO 算法^[17]实现对智能体的建模。PPO 算法鲁棒性高,调参简单,能够在连续动作空间中获得较为稳定的训练结果。PPO 算法通过对原始 Actor 结构做参考镜像以指导学习率大小,目标函数为

$$L^{\text{CLIP}}(\theta) = \widehat{E}_t \cdot \left\{ \min \left(r_t(\theta) \widehat{A}_t, \text{clip} [r_t(\theta), 1 - \epsilon, 1 + \epsilon] \widehat{A}_t \right) \right\} \quad (1)$$

式中: ϵ 为超参数; $r_t(\theta)$ 为新旧策略比; \widehat{A}_t 为新策略较旧策略的优势函数。

$$r_t(\theta) = \frac{\pi_{\theta_{\text{new}}}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \quad (2)$$

式中: $\pi_{\theta_{\text{new}}}(a_t|s_t)$ 为新策略; $\pi_{\theta_{\text{old}}}(a_t|s_t)$ 为旧策略; a_t 为 t 时刻的动作值; s_t 为 t 时刻的状态值。

由式(1)可以看出:PPO 算法在采集状态数据随机梯度优化的基础上增加了带有截断概率比的目标函数,可以避免策略参数大幅度更新导致的算法波动,具有更好的鲁棒性和效率。

在本文的应用场景中,Actor 的输入为态势特征信息,输出为飞机的行为决策;Critic 的输入为态势特征信息,输出为当前态势的价值,其通过和环境交互获得奖惩信息修正价值评估网络,并指导 Actor 做出正确的决策。算法流程如图 1 所示。

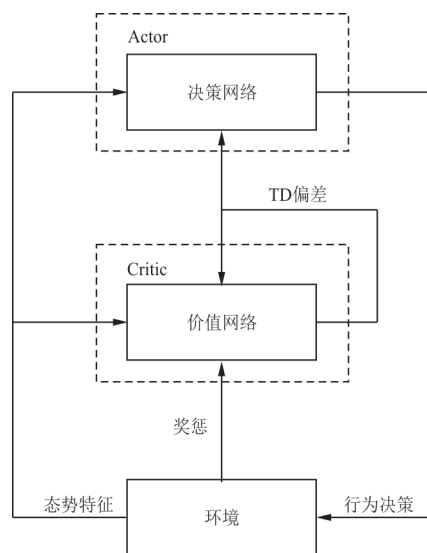


图 1 算法流程图

Fig. 1 Algorithm flow chart

本文的 Actor 和 Critic 均采用平均深度为 8 层的全连接结构,并使用层归一化(Layer Normalization, 简称 LN)和残差模块(Residual Net, 简称 Resnet)防止训练过程中的梯度爆炸或梯度消失现象,加快网络的收敛速度。

本文使用 5 个单独的网络分别作为 1 架有人机和 4 架无人机的决策中心,每个网络的输入为综合态势特征信息,输出为每个智能体的行为决策,奖惩函数独立。此种方式可以避免使用一个网络和一套奖惩机制进行决策时,某智能体做出错误决策,但整体评分上升导致的错误决策倾向及局部最优的问题。

1.2 状态空间描述

本文基于三维空间中 5v5 有人/无人机协同作战的想定环境,作战双方均为 1 架携带 4 枚导弹的有人机和 4 架携带 2 枚导弹的无人机,所有飞机初始部署均在空中,不考虑飞机的起飞和降落。将每架飞机和每枚导弹均简化为一个质点。飞机状态 s 能够通过三维坐标、姿态、速度和加速度等物理量描述,在任意时刻, s 可以由一个九元组表示,其向量形式为

$$s = [x, y, z, \theta, \psi, \phi, v, N_x, N_z] \quad (3)$$

式中: x, y, z 分别为飞机在三维空间坐标系下的坐标值, x 轴指向正东方向, y 轴指向正北方向, z 轴指向垂直向上方向; θ, ψ, ϕ 分别为飞机的俯仰角、偏航角和滚转角; v 为飞行速度; N_x 和 N_z 分别为飞机的切向和法向过载。

导弹状态能够通过三维坐标、姿态和速度等物理量描述,为了网络能够更好地学习导弹与目标的相对关系,将导弹状态 m 使用发射状态、目标、飞行时间、飞行距离等相对物理指标描述,在任意时刻, m 可以由一个六元素的十元组描述,其向量形式为

$$m = [S_{\text{tatus}}, T_{\text{arget}}, T_{\text{ime}}, D_{\text{flying}}, D_{\text{target}}, S_{\text{peed}}] \quad (4)$$

式中: S_{tatus} 为导弹状态, 1 表示挂载, 0 表示发射; T_{arget} 为 5 维特征向量形式表示的攻击目标, 分别对

应 5 架敌方飞机的索引, 值最大者将被攻击; T_{ime} 为导弹飞行时间; D_{flying} 为导弹飞行距离; D_{target} 为导弹相对目标距离; S_{peed} 为导弹方向与目标方向的投影, 数值越大代表导弹对目标的威胁越大。

有人机状态可由 1 架飞机状态和 4 枚导弹状态叠加, 形成 49 维的状态向量。无人机状态可由 1 架飞机状态和 2 枚导弹状态叠加, 形成 29 维的状态向量。敌我双方可以分别用 165 维的状态向量表示。为了便于网络统计战局中有利或不利态势的积累情况, 提供 1 维全局的时间状态输入。因此, 整个空战系统的状态向量共 331 维, 在输入给网络计算前, 均根据各自的极值归一化到 $[-1, 1]$ 区间。

1.3 态势评估奖惩函数

态势评估奖惩函数以实时态势为输入, 以对抗的胜负条件为基础评价每一步决策的好坏, 对于能促进我方进入有利态势的决策予以奖励, 否则予以惩罚。对抗的胜负规则如表 1 所示。

表 1 胜负规则
Table 1 Success and failure rules

序号	胜利条件	失败条件
1	敌方有人机战损	我方有人机战损
2	敌方无弹	我方无弹
3	到达终止时间, 我方剩余无人机数量大于敌方	到达终止时间, 我方剩余无人机数量小于敌方
4	到达终止时间, 双方剩余无人机与导弹数相同, 我方控制作战中心区域时间大于敌方	到达终止时间, 双方剩余无人机与导弹数相同, 我方控制作战中心区域时间小于敌方

本文设置的奖惩函数将引导飞机向避免战损并且尽可能击落敌有人机的方向做出决策。奖惩的优先级按飞机类型分: 有人机 > 无人机; 按状态和行为分: 存活/战损 > 攻击/被攻击 > 制导/占领中心区域 > 位置引导/逃脱, 具体规则和数值设置如表 2 所示。

表 2 奖惩规则
Table 2 Reward and punishment rules

奖惩类型	条件	有人机奖惩	无人机奖惩
存活/损毁	自身战损	-5	-1
	击毁敌有人机的导弹发射者	+5	+5
	击毁敌无人机的导弹发射者	+1	+1
攻击/被攻击	向敌有人机发射导弹	+0.5	+0.5
	向敌无人机发射导弹	+0.1	+0.1
	被敌方导弹瞄准,飞机姿态朝向与导弹朝向同向扣分,偏离加分	见式(5)	见式(5)
制导	为我方已发射导弹制导,制导目标为敌有人机	+0.02	+0.02
	为我方已发射导弹制导,制导目标为敌无人机	+0.01	+0.01
中心区域	进入中心区域	+0.01,累积加分截断为[0,1]	+0.01,累积加分截断为[0,1]
	引导无人机打击/迷惑进入我方有人机威胁的敌方飞机	-	见式(6)
位置引导	没有上述任何奖惩时,引导飞机飞向敌有人机,飞机朝向与其和敌有人机连线夹角小于30°时	+0.001	+0.001

奖惩数值的量级由优先级和行为特点共同决定,如:高优先级、一次性的存活状态,数值量级最高;可积分累加奖惩的行为会根据累加方式(次数或时间)决定数值量级。此种设置方式能够保证完全积分累加后的奖惩数值与最高量级相同,并按优先级排序,最终实现在引导中间行为向有利态势发展的同时,防止发生由中间行为过度主导训练过程,影响最终结果的问题。

为了使决策更加灵活,实现有/无人机的协同作战,加入被攻击时的脱离奖励、导弹制导奖励和基于敌机对我方有人机威胁度的无人机目标分配和打击/诱导奖励。

在飞机被敌方导弹攻击时,如果飞机处于可逃逸区,则判断飞机姿态朝向与导弹朝向的关系:同向扣分,偏离加分。奖惩规则为

$$R_{\text{rewards}} = 0.005 \times \frac{90 - \theta}{90} \quad (5)$$

式中: θ 为飞机与导弹之间的脱离角,即飞机朝向与导弹与飞机连线的夹角。

飞机与导弹的相对关系示意图如图2所示。当我方飞机发射导弹进行攻击,且导弹处于中制导状态时,我方任何飞机若对该导弹进行制导,即保证目标在飞机火控雷达的视场范围和距离范围内,则可获得制导奖励。制导奖励低于完全偏移下的脱离奖励,能够引导网络习得接力制导的作战模式。

当敌机进入对我方有人机有威胁的距离范围时,根据其雷达范围与携弹数量进行威胁度排序,范围越大、携弹数量越高则威胁度越高;同时对我

方无人机根据携弹数量进行能力排序,携弹数量越高能力越高,使用高能力的我方战机对抗高威胁敌机,当我方飞机能力值相同时优先派遣距离最近单位,并根据被派遣无人机与目标敌机的距离缩减值对其进行奖励,以实现对我方有人机的掩护。奖励规则为

$$R_{\text{rewards}} = 0.005 \times [D(t) - D(t-1)] \quad (6)$$

式中: $D(t)$ 为第 t 秒时,被派遣无人机与目标敌机的直线距离。

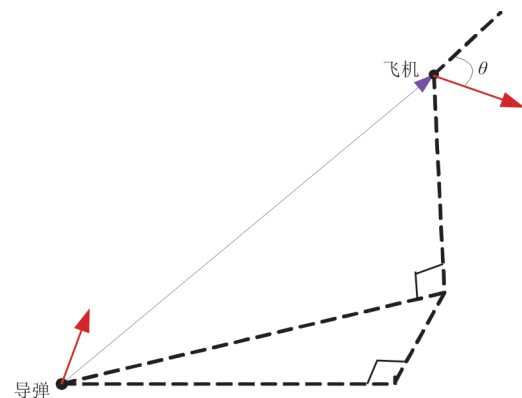


图 2 飞机与导弹相对关系示意图
Fig. 2 Schematic diagram of relative relationship between aircraft and missile

1.4 机动动作决策

将飞机的战斗决策分为机动类和攻击类两类。机动类决策动作使用由坐标、速度和过载定义的六元组,完成航向改变、高度改变与稳定飞行,其向量形式为

$$A_{\text{ction}} = [x, y, z, v, N_x, N_z] \quad (7)$$

式中: x, y, z 分别为飞机拟到达的三维空间坐标值; v 为拟采用的飞行速度; N_x 和 N_z 分别为飞机拟

采用的切向和法向过载。

飞机在三维空间中以 Δt 时间间隔更新的运动微分方程为

$$\begin{cases} v(t+1) = v(t) + \Delta t \cdot g [N_x(t) - \sin \theta(t)] \\ \phi(t+1) = \phi(t) + \Delta t \cdot N_z(t) \cdot \frac{g}{v(t)} \cdot \frac{\sin \phi(t)}{\cos \theta(t)} \\ \theta(t+1) = \theta(t) + \Delta t \cdot \frac{g}{v(t)} \cdot [N_z \cos \phi(t) - \cos \theta(t)] \\ x(t+1) = x(t) + \Delta t \cdot v(t) \cos \theta(t) \cos \phi(t) \\ y(t+1) = y(t) + \Delta t \cdot v(t) \cos \theta(t) \sin \phi(t) \\ z(t+1) = z(t) + \Delta t \cdot \sin \theta(t) \end{cases} \quad (8)$$

攻击类决策是由攻击指令和攻击目标组成的六元组,其向量形式为

$$A_{\text{ttack}} = [L_{\text{aunch}}, T_{\text{arget}}] \quad (9)$$

式中: L_{aunch} 为由 0/1 组成的发射指令; T_{arget} 为 one hot 形式的 5 元组, 对应敌方飞机的 ID 索引。

每架飞机的网络输出均为经过归一化处理、值在 $[-1, 1]$ 之间的 12 维向量, 在由网络输出封装指令集时需根据各参数的实际范围做伸缩变换。

本文在网络训练过程中, 设置训练参数 batch size 为 64, Actor 学习率为 0.2×10^{-4} , Critic 学习率为 0.4×10^{-4} 。模型训练过程中奖励震荡收敛, 为便于查看趋势, 使用滑动平均法 (滑动窗为 10) 展示收敛曲线, 网络经 800 场训练后的奖励收敛结果如图 4 所示。

2 仿真与分析

仿真实验的硬件环境为: Intel i9-11900K@3.5 GHz CPU, 64 G 内存, NVIDIA Geforce RTX 3090 显卡, 24 G 显存; 软件环境为: Ubuntu 20.04 LTS 操作系统, Python 3.7, Tensorflow 1.14.0。智能体仿真训练框架如图 3 所示。

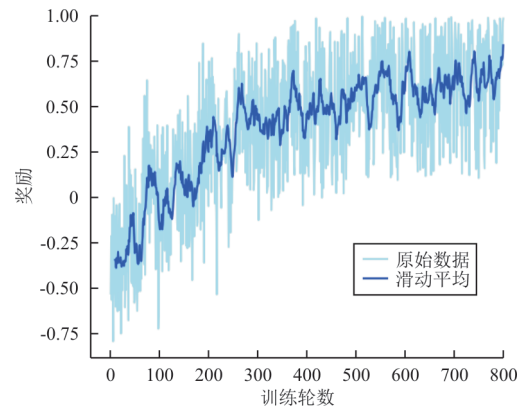


图 4 奖励收敛结果

Fig. 4 Reward convergence results

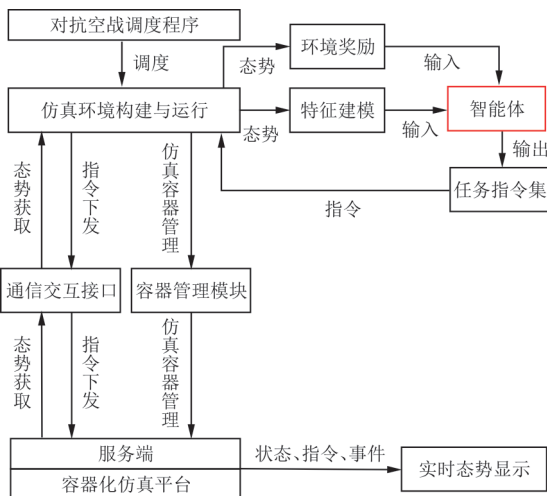


图 3 智能体仿真训练框架

Fig. 3 Agent simulation training framework

使用规则智能体与强化学习智能体进行对抗训练。规则智能体以击落敌有人机为主要原则。初始时, 四架无人机以有人机为中心, 等间距排布在其两侧。为增加敌方的应对难度, 每场有人机的初始位置随机。战斗开始的前 300 s, 有人机在原地盘旋, 无人机向高处及敌有人机所在方向飞行, 之后, 有人机和无人机均朝向敌有人机飞行。在飞行过程中, 如果有敌机进入攻击范围, 则发射导弹攻击。为保证导弹使用效率, 攻击前会对空中正在引导攻击的导弹和余弹量进行统计。当正在引导攻击敌有人机的导弹数量大于等于 4, 攻击敌无人机的导弹数量大于等于 2 时, 不再对该敌机

发射导弹。当全机队的余弹量只剩 1 枚时,判断我方是否处于能够获胜的有利态势,若态势被动,则追随敌有人机并发射最后一枚导弹。

在仿真训练实验中,仿真时间步长 $S_{\text{step}}=1\text{ s}$,即每隔 1 s 敌我双方获取一次当前态势,更新决策。作战范围为 $300\text{ km}\times 300\text{ km}$,战局时长为 20 min,根据表 1 中的规则判定胜负。

由于本文使用五个独立的网络进行训练,训练过程中会出现个别网络收敛快、个别网络难收敛的情况。针对这种情况,需要固定已收敛的网络参数,对未收敛的网络单独训练。当五个网络均达到较好效果后再一起训练至全部收敛。

收敛后的强化学习智能体对规则智能体的 100 局对抗结果如图 5~图 6 所示,其中图 5 为智能体扮演红方的对抗结果,图 6 为智能体扮演蓝方的对抗结果。

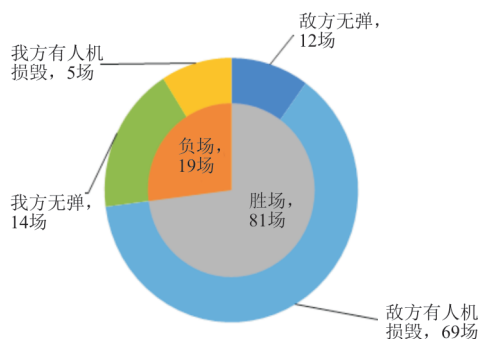


图 5 红方智能体对抗结果

Fig. 5 Red agent confrontation results

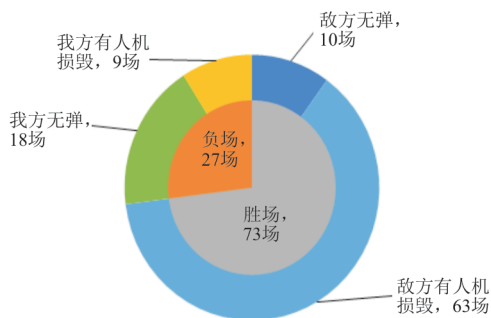


图 6 蓝方智能体对抗结果

Fig. 6 Blue agent confrontation results

从图 5~图 6 可以看出:强化学习智能体的胜率能够达到 70%~80%,对抗结束条件集中在有人机损毁和无弹两种模式,并无到达截止时间需进行二级判断或平局的情况;在我方胜利的对战场次中,敌有人机损毁占 85% 左右,敌方无弹占

15% 左右;在我方失败的对战场次中,我方有人机损毁占 30% 左右,我方无弹占 70% 左右,可见本文的强化学习智能体能够按照奖励规则的引导进行学习,向避免我方战损和倾向击落敌有人机的方向做出决策。面对环境中不同位置进攻的敌方,本文的智能体能够习得泛化的对敌策略,并达到最终击毁敌有人机的目标。

3 结 论

本文对 5v5 有/无人机协同空战决策进行研究,提出基于近端策略优化的学习方法,并与基于规则的智能体进行对战仿真训练。结果表明,基于本文方法训练得到的智能体,其行为决策具有实时性和泛化性,能够根据奖惩函数的设置及一定阶段的训练,实现对敌方的主动攻击、制导,面对威胁时的规避或诱敌等,最终形成一套鲁棒的空战决策方法。

本文提出的算法虽然在 5v5 有/无人机协同空战决策中取得了一定的效果,但实际的空战态势环境会更加复杂,未来可以结合专家经验重塑奖惩函数,以实现更加复杂的战术设计与行为设计,提升智能体的决策能力。

参 考 文 献

- [1] WANG R, GAO Z. Research on decision system in air combat simulation using maneuver library [J]. Flight Dynamic, 2009, 27(6): 72-75.
- [2] 傅莉, 王晓光. 无人战机近距空战微分对策建模研究 [J]. 兵工学报, 2012, 33(10): 1210-1216.
FU Li, WANG Xiaoguang. Research on close air combat modeling of differential games for unmanned combat air vehicles [J]. Acta Armamentarii, 2012, 33(10): 1210-1216. (in Chinese)
- [3] 王宁宇, 苏山, 崔乃刚, 等. 基于微分对策的多飞行器协同最优分配方法 [J]. 战术导弹技术, 2021(6): 130-138.
WANG Ningyu, SU Shan, CUI Naigang, et al. Multi-aircraft cooperative optimal allocation method based on differential games [J]. Tactical Missile Technology, 2021(6): 130-138. (in Chinese)
- [4] HA J S, CHAE H J, CHOI H L. A stochastic game-theoretic approach for analysis of multiple cooperative air combat [C]// 2015 American Control Conference (ACC). Chicago, IL, USA: IEEE, 2015: 3728-3733.
- [5] 钱炜祺, 车竞, 何开锋. 基于矩阵博弈的空战决策方法

- [C]// 第二届中国指挥控制大会. 北京: 中国指挥与控制学会, 2014: 408-412.
- QIAN Weiqi, CHE Jing, HE Kaifeng. Air combat decision method based on game-matrix approach[C]// The Second China Command and Control Conference. Beijing: China Command and Control Society, 2014: 408-412. (in Chinese)
- [6] 况立群, 李思远, 冯利, 等. 深度强化学习算法在智能军事决策中的应用[J]. 计算机工程与应用, 2021, 57(20): 271-278.
- KUANG Liqun, LI Siyuan, FENG Li, et al. Application of deep reinforcement learning algorithm on intelligent military decision system [J]. Computer Engineering and Applications, 2021, 57(20): 271-278. (in Chinese)
- [7] 单圣哲, 杨孟超, 张伟伟, 等. 自主空战连续决策方法[J]. 航空工程进展, 2022, 13(5): 47-58.
- SHAN Shengzhe, YANG Mengchao, ZHANG Weiwei, et al. Continuous decision-making method for autonomous air combat[J]. Advances in Aeronautical Science and Engineering, 2022, 13(5): 47-58. (in Chinese)
- [8] 左家亮, 杨任农, 张滢, 等. 基于启发式强化学习的空战机动智能决策[J]. 航空学报, 2017, 38(10): 212-225.
- ZUO Jialiang, YANG Rennong, ZHANG Ying, et al. Intelligent decision-making in air combat maneuvering based on heuristic reinforcement learning[J]. Acta Aeronautica et Astronautica Sinica, 2017, 38(10): 212-225. (in Chinese)
- [9] ZHANG X, LIU G, YANG C, et al. Research on air confrontation maneuver decision-making method based on reinforcement learning[J]. Electronics, 2018, 7(11): 279-298.
- [10] ERNEST N, COHEN K, KIVELEVITCH E, et al. Genetic fuzzy trees and their application towards autonomous training and control of a squadron of unmanned combat aerial vehicles[J]. Unmanned Systems, 2015, 3(3): 185-204.
- [11] TOUBMAN A, ROESSINGH J J, SPRONCH P, et al. Transfer learning of air combat behavior[C]// Proceedings of the IEEE International Conference on Machine Learning and Applications. Miami, US: IEEE, 2015: 226-231.
- [12] 吴宜珈, 赖俊, 陈希亮, 等. 强化学习算法在超视距空战辅助决策上的应用研究[J]. 航空兵器, 2021, 28(2): 55-61.
- WU Yijia, LAI Jun, CHEN Xiliang, et al. Research on the application of reinforcement learning algorithm in decision support of beyond-visual-range air combat[J]. Aero Weaponry, 2021, 28(2): 55-61. (in Chinese)
- [13] HINTON G E, OSINDERO S, TEH Y W. A fast learning algorithm for deep belief nets[J]. Neural Computation, 2006, 18(7): 1527-1554.
- [14] SILVER D, HUANG A, MADDISON C J, et al. Mastering the game of go with deep neural networks and tree search [J]. Nature, 2016, 529: 484-489.
- [15] SILVER D, LEVER G, HEESS N, et al. Deterministic policy gradient algorithms[C]// The 31st International Conference on Machine Learning. Beijing: IEEE, 2014: 387-395.
- [16] MNIH V, BADIA A P, MIRZA M, et al. Asynchronous methods for deep reinforcement learning[C]// The 33rd International Conference on Machine Learning. New York, US: IEEE, 2016: 1928-1937.
- [17] SCHULMAN J, WOLSKI F, DHARIWAL P, et al. Proximal policy optimization algorithms [EB/OL]. (2017-07-20)[2022-06-11]. <https://arxiv.org/abs/1707.06347>.

作者简介:

张博超(1982-),男,硕士,高级工程师。主要研究方向:机载软件算法设计,嵌入式软件研发,软件工程化管理等。

温晓玲(1979-),女,硕士,高级工程师。主要研究方向:嵌入式软件研发,集成验证和型号软件技术状态管理等。

刘 璐(1995-),女,硕士,工程师。主要研究方向:机载软件算法设计,嵌入式软件研发,软件工程化管理等。

张雅茜(1978-),女,硕士,高级工程师。主要研究方向:航空嵌入式软件研发,集成验证和型号软件技术状态管理。

王宏光(1976-),女,学士,高级工程师。主要研究方向:机载软件算法设计,嵌入式软件研发,软件工程化管理等。

(编辑:马文静)