

文章编号: 2097-1974(2025)06-0001-09

DOI: 10.7654/j.issn.2097-1974.20250601

基于改进型TD3强化学习的高速飞行器姿态控制

王伟丽, 黄万伟, 刘晓东, 路坤锋, 贾晨辉

(北京航天自动控制研究所, 宇航智能控制技术全国重点实验室, 北京, 100854)

摘要: 针对高速飞行器再入段面临的强非线性、高不确定性以及参数快时变等挑战, 结合航天器智能化发展需求, 提出了一种改进型的双延迟深度确定性策略梯度 (Twin Delayed Deep Deterministic Policy Gradient, TD3) 端到端智能姿态控制方法。为解决TD3算法在姿态控制学习过程中存在训练不稳定、收敛困难的问题, 在其马尔可夫决策过程中, 设计了混合奖励机制, 融合连续跟踪误差惩罚和稀疏任务完成奖励, 协同引导智能体收敛; 在其训练过程中, 引入基于现代控制理论的先验知识约束, 提出了基于行为克隆的Actor网络优化更新策略, 以平衡专家经验模仿与累计回报最大化目标。仿真结果表明, 在14种参数偏差组合的工况下, 所提方法能够精确跟踪三通道姿态指令。

关键词: 高速飞行器; 姿态控制; 深度强化学习; 行为克隆; 强适应控制

中图分类号: V448.2

文献标识码: A

Attitude Control of High-speed Vehicles Based on Improved TD3 Reinforcement Learning

WANG Weili, HUANG Wanwei, LIU Xiaodong, LU Kunfeng, JIA Chenhui

(National Key Laboratory of Science and Technology on Aerospace Intelligent Control, Beijing Aerospace Automatic Control Institute, Beijing, 100854)

Abstract: To address the challenges of strong nonlinearity, high uncertainty, and rapid time-varying parameters during the reentry phase of high-speed vehicles, this study proposes an end-to-end intelligent attitude control method based on an improved Twin Delayed Deep Deterministic Policy Gradient algorithm, aligned with the demands of intelligent spacecraft development. To overcome the issues of training instability and convergence difficulties in TD3-based attitude control learning, two key innovations are introduced: a hybrid reward mechanism combining continuous tracking error penalties and sparse task-completion rewards is designed within the Markov Decision Process framework to synergistically guide agent convergence. Prior knowledge constraints derived from modern control theory are incorporated into the training process, proposing a behavior cloning-based optimization strategy for the Actor network to balance expert experience imitation and cumulative reward maximization. Simulation results show that the proposed method can accurately track the three-channel attitude commands under 14 combinations of parameter deviations.

Keywords: high-speed vehicles; attitude control; deep reinforcement learning; behavior cloning; strongly adaptive control

0 引言

高速飞行器凭借其全空域机动、宽速域巡航和强突防能力的优势, 已成为现代远程精确打击体系的核心装备^[1]。然而, 其飞行包线内表现出的强非线性动力学特性、多通道耦合效应以及由气动热或结构形变引发的快时变参数, 使得传统基于精确数学模型的控制方法设计面临理论局限^[2]。尤其在再入段高动态

环境下, 飞行器同时承受极端气动载荷、复杂干扰和模型不确定性, 基于先验知识的经典控制方法(如增益调度PID、鲁棒自适应控制)进而设计可靠的姿态控制器变得困难。

为应对这些挑战, 并响应航天器智能化发展趋势^[3], 智能飞行控制(Intelligent Flight Control, IFC)技术应运而生。其中, 数据驱动方法因其对模

型依赖程度低的优势,正逐渐成为IFC领域的研究重点。深度强化学习(Deep Reinforcement Learning, DRL)因其特有的“环境交互-自主优化”机制,展现出解决复杂控制问题的独特潜力^[4-6]。目前DRL在高速飞行器控制中的应用主要呈现3个研究方向^[7]:控制参数自适应整定^[8-9]、不确定性补偿控制^[10]以及端到端自主控制^[11],形成了该领域新的技术突破点。

基于DRL的端到端控制架构通过直接从原始传感器数据学习控制策略以生成执行器命令,无需人工特征提取或控制律设计,实现高速飞行器主控制律的高自主设计,有效降低对其地面设计模型的依赖程度。然而,DRL在高速飞行器再入段三通道姿态控制中的应用研究仍处于探索阶段,其核心挑战源于两大特性:一是再入过程中气动参数与惯性参数的快时变特性导致系统动力学环境剧烈变化;二是三通道间的强耦合效应使得控制策略的训练难以稳定收敛。针对这些问题,现有研究主要从两个方向突破:在算法层面,学者们通过优化网络架构设计^[12]等方式和改进经验回放机制^[13]来提升收敛性能;在训练策略层面,结合行为克隆(Behavior Cloning, BC)技术以修正策略偏差^[14]。

综上所述,本研究针对高速飞行器再入段端到端姿态控制问题,创新性提出了基于知识引导的双延迟深度确定性策略梯度(Knowledge-Guided Twin Delayed Deep Deterministic Policy Gradient, KG-TD3)算法,该算法通过融合现代控制理论中的先验知识与双延迟深度确定性策略梯度(Twin Delayed Deep Deterministic Policy Gradient, TD3)算法的数据驱动特性,构建了新的混合驱动控制架构。这种知识嵌入式的DRL范式不仅能够解决纯数据驱动方法在复杂动态环境中的训练不稳定问题,同时保留了端到端控制的自适应优势,为高速飞行器的智能控制提供了新的技术途径。

1 模型建立

本文研究对象为升力式面对称无动力高速飞行器,其采用BTT控制模式,控制物理量为攻角 α 、侧滑角 β 和速度倾侧角 γ_v 。典型的无动力高速飞行器有美国的HTV-2、SR-72等,如图1所示。

根据升力式面对称高速飞行器自身及其再入段飞

行环境特点,提出如下可行性假设条件:忽略地球自转的影响,此时发射惯性坐标系与地面坐标系始终重合,而且不需考虑离心惯性和哥氏惯性的作用;将地球视为均质圆球,忽略地球扁率以及切向引力加速度的影响;惯量积 $J_{x,y,z}$ 为小量,且忽略不计;将飞行器视为刚体,即不考虑弹性影响。



HTV-2



SR-72

图1 典型的无动力高速飞行器

Fig.1 Typical unpowered high-speed vehicle

参考文献[15]建立了飞行器六自由度数学模型,并写为如式(1)所示的仿射非线性形式,以便于控制系统设计。

$$\begin{cases} \dot{\mathbf{x}}_\Omega = \mathbf{F}_\Omega + \mathbf{G}_\Omega \mathbf{x}_\omega + \mathbf{d}_\Omega \\ \dot{\mathbf{x}}_\omega = \mathbf{F}_\omega + \mathbf{G}_\omega \mathbf{u} + \mathbf{d}_\omega \end{cases} \quad (1)$$

式中 \mathbf{x}_Ω 为姿态环状态量,即攻角 α 、侧滑角 β 、速度倾侧角 γ_v ; \mathbf{x}_ω 为姿态角速度环状态量,即滚转角速度 ω_x 、偏航角速度 ω_y 、俯仰角速度 ω_z ; \mathbf{u} 表示控制输入,即滚转舵偏角 δ_x 、偏航舵偏角 δ_y 、俯仰舵偏角 δ_z ; \mathbf{F}_Ω , \mathbf{F}_ω 表示受控对象的集中动力学; \mathbf{G}_Ω , \mathbf{G}_ω 为控制信号的增益,描述了控制信号对系统动力学的影响; $\mathbf{d}_\Omega = [d_\alpha \ d_\beta \ d_{\gamma_v}]^T$, $\mathbf{d}_\omega = [d_{\omega_x} \ d_{\omega_y} \ d_{\omega_z}]^T$ 为外界干扰。

在实际工程中,由于高速飞行器飞行过程中存在参数摄动大、外界干扰严重等问题,则将飞行器模型写为如式(2)所示的仿射非线性系统:

$$\begin{cases} \dot{\mathbf{x}}_\Omega = \bar{\mathbf{F}}_\Omega + \bar{\mathbf{G}}_\Omega \mathbf{x}_\omega + \zeta_\Omega \\ \dot{\mathbf{x}}_\omega = \bar{\mathbf{F}}_\omega + \bar{\mathbf{G}}_\omega \mathbf{u} + \zeta_\omega \end{cases} \quad (2)$$

式中 $\mathbf{F}_\Omega = \bar{\mathbf{F}}_\Omega + \Delta\mathbf{F}_\Omega$, $\mathbf{F}_\omega = \bar{\mathbf{F}}_\omega + \Delta\mathbf{F}_\omega$, $\mathbf{G}_\Omega = \bar{\mathbf{G}}_\Omega + \Delta\mathbf{G}_\Omega$,

$G_\omega = \bar{G}_\omega + \Delta G_\omega$ 。 \bar{F}_Ω , \bar{G}_Ω 是姿态环标称状态下的已知模型, \bar{F}_ω , \bar{G}_ω 是姿态角速度环标称状态下的已知模型, 具体如式 (3) 和式 (4) 所示。 ΔF_Ω , ΔG_Ω , ΔF_ω , ΔG_ω 为未知动态。 ζ_Ω 和 ζ_ω 表示模型未知动态、参数摄动、外界干扰等系统未知项, $\zeta_\Omega = \Delta F_\Omega + \Delta G_\Omega x_\omega + d_\Omega$, $\zeta_\omega = \Delta F_\omega + \Delta G_\omega u + d_\omega$, 该项的存在是导致系统性能下降的主要原因, 需采用自适应强抗扰的控制器应对, 本文采用的是KG-TD3智能控制器。

$$\begin{cases} \dot{x}_\Omega = [\dot{\alpha} \quad \dot{\beta} \quad \dot{\gamma}_v]^\top \\ x_\omega = [\omega_x \quad \omega_y \quad \omega_z]^\top \\ \bar{F}_\Omega = \begin{bmatrix} \frac{L - mg \cos \theta \cos \gamma_v}{mV \cos \beta} \\ \frac{Z + mg \cos \theta \sin \gamma_v}{mV} \\ \frac{[L(\tan \beta + \tan \theta \sin \gamma_v) + Z \tan \theta \cos \gamma_v - mg \cos \theta \tan \beta \cos \gamma_v]}{mV} \end{bmatrix} \\ \bar{G}_\Omega = \begin{bmatrix} -\cos \alpha \tan \beta & \sin \alpha \tan \beta & 1 \\ \sin \alpha & \cos \alpha & 0 \\ \cos \alpha \sec \beta & -\sin \alpha \sec \beta & 0 \end{bmatrix} \end{cases} \quad (3)$$

$$\begin{cases} \dot{x}_\omega = [\dot{\omega}_x \quad \dot{\omega}_y \quad \dot{\omega}_z]^\top \\ u = [\delta_x \quad \delta_y \quad \delta_z]^\top \\ \bar{F}_\omega = \begin{bmatrix} \frac{J_y - J_z}{J_x} \omega_y \omega_z + \frac{(C_{mx}^\alpha \alpha + C_{mx}^\beta \beta) qSl}{J_x} \\ \frac{J_z - J_x}{J_y} \omega_x \omega_z + \frac{C_{my}^\beta \beta qSl}{J_y} \\ \frac{J_x - J_y}{J_z} \omega_x \omega_y + \frac{C_{mz}^\alpha \alpha qSl}{J_z} \end{bmatrix} \\ \bar{G}_\omega = qSl \begin{bmatrix} \frac{1}{J_x} & 0 & 0 \\ 0 & \frac{1}{J_y} & 0 \\ 0 & 0 & \frac{1}{J_z} \end{bmatrix} \begin{bmatrix} C_{mx}^{\delta_x} & C_{mx}^{\delta_y} & C_{mx}^{\delta_z} \\ C_{my}^{\delta_x} & C_{my}^{\delta_y} & C_{my}^{\delta_z} \\ C_{mz}^{\delta_x} & C_{mz}^{\delta_y} & C_{mz}^{\delta_z} \end{bmatrix} \end{cases} \quad (4)$$

式中 m , V , θ 分别为飞行器的质量、速度和弹道倾角; L , Z 分别为气动升力和气动侧向力; g 为重力加速度; M_{x1} , M_{y1} 和 M_{z1} 分别为气动滚转力矩、偏航力矩和俯仰力矩; J_{x1} , J_{y1} 和 J_{z1} 为飞行器的主转动惯量; q , S , l 分别为动压、气动参考面积和参考长度; $C_{mx}^{(\cdot)}$, $C_{my}^{(\cdot)}$, $C_{mz}^{(\cdot)}$ 分别为相对于 (\cdot) 的滚动、偏航和俯仰力矩系数。

至此, 面向控制的高速飞行器三通道姿态运动数学模型构建完成。接下来, 将根据该模型研究知识与数据融合的智能姿态控制方法。

2 知识与数据融合的智能控制器设计

在深度强化学习中, TD3 算法虽然在连续控制任务中表现出色, 但仍面临探索效率低、训练初期不稳定以及局部最优陷阱等问题。为此, 本文引入基于知识的控制器约束, 即动态面控制器 (Dynamic Surface Control, DSC) (本文将其定义为“专家控制器”), 并结合行为克隆方法, 提出了基于知识引导的TD3 (KG-TD3) 算法。其中, “知识”与“数据”均依据文献 [16] 界定。

2.1 基于知识的控制器设计

动态面控制器依赖于精确、解析的飞行器数学模型, 故依据文献 [16] 可称为基于知识的控制器。

对于非线性系统, 定义跟踪误差 s_Ω 和其微分:

$$\begin{cases} s_\Omega = x_\Omega - x_{\Omega d} \\ \dot{s}_\Omega = \bar{F}_\Omega + \bar{G}_\Omega x_\omega + \zeta_\Omega - \dot{x}_{\Omega d} \end{cases} \quad (5)$$

式中 x_Ω 为系统实际状态; $x_{\Omega d}$ 为系统制导指令状态。

取虚拟控制输入 x_v :

$$x_v = -\bar{G}_\Omega^{-1}(\bar{F}_\Omega + W_\Omega s_\Omega - \dot{x}_{\Omega d} + \zeta_\Omega) \quad (6)$$

经一阶滤波器, 得到:

$$\tau \dot{x}_\tau + x_\tau = x_v \quad (7)$$

式中 τ 为待设计的滤波器系数; x_τ 为经过滤波器后的控制输入。

定义跟踪误差 s_ω 并对其微分:

$$\begin{cases} s_\omega = x_\omega - x_\tau \\ \dot{s}_\omega = \bar{F}_\omega + \bar{G}_\omega u + \zeta_\omega - \dot{x}_\tau \end{cases} \quad (8)$$

从而设计最终的控制律 u :

$$u = -\bar{G}_\omega^{-1}[\bar{F}_\omega + W_\omega s_\omega - \dot{x}_\tau + \zeta_\omega] \quad (9)$$

综上, 动态面姿态控制律:

$$\begin{cases} s_\Omega = x_\Omega - x_{\Omega d} \\ x_v = -\bar{G}_\Omega^{-1}(\bar{F}_\Omega + W_\Omega s_\Omega - \dot{x}_{\Omega d} + \zeta_\Omega) \\ \tau \dot{x}_\tau + x_\tau = x_v \\ s_\omega = x_\omega - x_\tau \\ u = -\bar{G}_\omega^{-1}(\bar{F}_\omega + W_\omega s_\omega - \dot{x}_\tau + \zeta_\omega) \end{cases} \quad (10)$$

式中 W_Ω 和 W_ω 均为正定矩阵。

DSC 控制器所得到的 u , 即滚转舵偏角 δ_x 、偏航舵偏角 δ_y 、俯仰舵偏角 δ_z , 将其视为“专家动作”, 引导后续设计的智能控制器训练环节。

2.2 基于知识引导的TD3智能控制器设计

2.2.1 马尔可夫决策过程模型设计

在训练智能控制器前, 需要建立高速飞行器再入段飞行的马尔可夫决策过程 (Markov Decision Process, MDP), 以创建环境与智能体之间的联系。

a) 状态空间。

状态空间包括智能体可以从环境中收集到的有用信息。在本研究中,我们优先考虑高速飞行器的可观状态,形成如式(11)所示的状态空间。

$$\mathbf{s}_t = [\boldsymbol{\rho}_\Omega \mathbf{e}_\Omega, \boldsymbol{\rho}_\omega \mathbf{e}_\omega]^\top = [e_\alpha, e_\beta, e_\gamma, e_{\omega_x}, e_{\omega_y}, e_{\omega_z}]^\top \quad (11)$$

式中 \mathbf{e}_Ω 表示当前时刻姿态角跟踪误差; \mathbf{e}_ω 表示当前时刻姿态角速度跟踪误差; $\boldsymbol{\rho}_\Omega$ 和 $\boldsymbol{\rho}_\omega$ 为归一化正定对角矩阵,用于保证状态量的尺度大小基本相同。

b) 动作空间。

本研究为高速飞行器的三通道姿态控制,故智能体直接学习控制指令,设计如式(12)所示动作空间。

$$\mathbf{a}_t = [\delta_x, \delta_y, \delta_z]^\top \quad (12)$$

式中 δ_x 为滚转舵偏角; δ_y 为偏航舵偏角; δ_z 为俯仰舵偏角。

此外,考虑气动舵作动范围的物理限制,舵偏角的幅值应满足给定的约束范围。

$$-35^\circ \leq \delta_i < 35^\circ, i = x, y, z \quad (13)$$

c) 奖励函数。

奖励函数被设计为连续和稀疏奖惩的混合函数,其组成部分可表示如下:

1) 姿态角和姿态角速度误差惩罚。

$$P_1 = K_1(|e_\alpha| + |e_\beta| + |e_\gamma|) + K_2(|e_{\omega_x}| + |e_{\omega_y}| + |e_{\omega_z}|) \quad (14)$$

2) 姿态角误差奖励。

$$R_1 = K_3 e^{-\eta_1(|e_\alpha| + |e_\beta| + |e_\gamma|)} \quad (15)$$

3) 动作及动作变化率抖动惩罚。

首先,为抑制过大动作,对动作的绝对值之和进行惩罚;其次,若当前时刻的动作为 $\delta_{x,c}$, $\delta_{y,c}$, $\delta_{z,c}$, 上一时刻动作为 $\delta_{x,l}$, $\delta_{y,l}$, $\delta_{z,l}$, 动作变化率可依次写为 $\Delta\delta_x = \delta_{x,c} - \delta_{x,l}$, $\Delta\delta_y = \delta_{y,c} - \delta_{y,l}$, $\Delta\delta_z = \delta_{z,c} - \delta_{z,l}$, 惩罚相邻时间步动作的变化幅度,以抑制高频抖动。

$$P_2 = K_4 e^{-\eta_2(|\delta_x| + |\delta_y| + |\delta_z|)} + K_5 e^{-\eta_3(|\Delta\delta_x| + |\Delta\delta_y| + |\Delta\delta_z|)} \quad (16)$$

4) 动作安全性惩罚。

若动作超出阈值 δ_M , 则对超限部分进行二次惩罚。

$$P_3 = K_6 \sum_i (|\delta_i| - \delta_M)^2, \text{ if } \forall |\delta_i| > \delta_M (i = x, y, z) \quad (17)$$

5) 高精度跟踪奖励。

$$R_2 = K_7, \text{ if } \sum_i (|e_\alpha| + |e_\beta| + |e_\gamma|) < M \quad (18)$$

最终得到混合奖励函数:

$$R = -\sum_{j=1}^3 P_j + \sum_{r=1}^2 R_r \quad (19)$$

式中 $K_l \in \mathbb{R}^+$, $\forall l \in \{1, 2, \dots, 7\}$; $\eta_1, \eta_2, \eta_3 \in \mathbb{R}^+$ 均为奖惩系数; $M \in \mathbb{R}^+$ 为设计要求所提的姿态角值。

2.2.2 智能控制器训练

a) Actor网络更新策略的改进。

为解决TD3算法在复杂任务中Actor网络可能面临因探索不足或训练初期Critic不准确,从而学到次优策略,且若单纯模仿“专家控制器”又无法超越“专家”水平,为解决该问题,本文提出了一种基于行为克隆的约束优化方法,从而对Actor网络的更新策略进行改进。

该方法的核心思想是将“专家控制器”的先验知识以软约束的形式融入DRL框架,具体来讲,在Actor网络策略优化目标中引入“专家动作”的行为克隆损失作为正则项,构建如式(20)所示的复合目标函数,以确保Actor输出的动作不会偏离专家动作太远,同时最大化Critic评估的 q 值,达到平衡“模仿专家”和“最大化累计回报”两个目标。

$$L(\theta) = \underbrace{-\mathbb{E}_{s_j \sim \mathcal{D}} [q(s_j, \hat{\mathbf{a}}_j; \boldsymbol{\omega}_{1,o})]}_{\text{强化学习目标项}} + \lambda \underbrace{\mathbb{E}_{s_j \sim \mathcal{D}} [\|\hat{\mathbf{a}}_j - \mathbf{a}_{j,k}\|^2]}_{\text{行为克隆约束项}} \quad (20)$$

式中 $L(\theta)$ 为损失函数,优化的目标即为最小化损失函数; $\mathbb{E}_{s_j \sim \mathcal{D}} [q(s_j, \hat{\mathbf{a}}_j; \boldsymbol{\omega}_{1,o})]$ 表示当前策略 $\boldsymbol{\omega}_{1,o}$ 在状态 s_j 下生成的动作 $\hat{\mathbf{a}}_j$, 由Critic网络评估得到 q 值后,在数据分布 \mathcal{D} (从经验回放缓冲区 \mathfrak{R} 中采样得到) 上的平均 q 值; $\mathbb{E}_{s_j \sim \mathcal{D}} [\|\hat{\mathbf{a}}_j - \mathbf{a}_{j,k}\|^2]$ 表示当前策略 $\boldsymbol{\omega}_{1,o}$ 在状态 s_j 下生成的动作 $\hat{\mathbf{a}}_j$ 与当前时刻的“专家动作” $\mathbf{a}_{j,k}$ 的均方误差在数据分布 \mathcal{D} 上的期望; λ 为行为克隆的权重因子。

b) KG-TD3算法训练框架。

KG-TD3算法训练与TD3有两点不同:一是提出了一种基于行为克隆的约束优化方法,利用式(20)进行Actor网络更新策略进行改进;二是为将智能控制器与“专家控制器”在时间尺度上对齐,需额外存储“专家动作”,本研究选取对经验回放缓冲区进行改进。具体来讲,在智能体与环境交互过程后,将智能体轨迹整理为 $(s_t, \mathbf{a}_t, r_t, s_{t+1}, \mathbf{a}_{t,k})$ 的五元组,即 t 时刻的状态 s_t 、智能体动作 \mathbf{a}_t 、奖励 r_t 以及 $t+1$ 时刻的状态 s_{t+1} 、“专家控制器”动作 $\mathbf{a}_{t,k}$ 。而后通过经验回放缓存区采用时,即可匹配同一时间的“专家动作”。

除上述两点外的KG-TD3训练环节与TD3别无二致,此处不再赘述,可参考文献[12]。KG-TD3完整的算法伪代码如下所示。

- 1: 随机初始化策略网络 $\mu(s; \theta)$ 和价值网络 $q_1(s, a; \omega_1), q_2(s, a; \omega_2)$
- 2: 初始化目标策略网络参数 $\theta^- \leftarrow \theta$ 和目标价值网络 $\omega_1^- \leftarrow \omega_1, \omega_2^- \leftarrow \omega_2$
- 3: 初始化经验回放缓冲区 \mathfrak{R}
- 4: **for** episode = 1 to M **do**
- 5: 初始化噪声 ϵ 用于噪声探索
- 6: 策略网络接收初始观察状态 s
- 7: **for** $t = 1$ to T **do**
- 8: 根据当前策略和探索噪声选择动作 $a_t = \mu(s_t; \theta) + \epsilon$
- 9: 执行动作 a_t , 获取奖励 r_t , 观测新状态 s_{t+1}
- 10: 经验存储: 在 \mathfrak{R} 中存储 $(s_t, a_t, r_t, s_{t+1}, a_{t,k})$
- 11: 经验回放: 从 \mathfrak{R} 中随机取出 N 个 $(s_j, a_j, r_j, s_{j+1}, a_{j,k})$
- 12: $\hat{a}_{j+1} = \mu(s_{j+1}; \theta_o^-) + \xi$
- 13: 两个目标价值网络预测: $\hat{q}_{i,j+1} = q(s_{j+1}, \hat{a}_{j+1}; \omega_{i,o}^-) (i = 1, 2)$
- 14: $\hat{y}_j = r_j + \gamma \cdot \min\{\hat{q}_{1,j+1}, \hat{q}_{2,j+1}\}$
- 15: 两个价值网络预测: $\hat{q}_{i,j} = q(s_j, a_j; \omega_{i,o})$
- 16: 最小化损失函数: $L(\omega_i) = \frac{1}{N} \left[(\hat{q}_{1,j} - \hat{y}_j)^2 + (\hat{q}_{2,j} - \hat{y}_j)^2 \right]$
- 17: 更新价值网络 $\omega_{i,e} \leftarrow \omega_{i,o}$
- 18: **if** $t \bmod k$ **then**
- 19: $\hat{a}_j = \mu(s_j; \theta_o)$
- 20: 最小化损失函数 $L(\theta)$, 更新策略网络 $\theta_e \leftarrow \theta_o$
- 21: 软更新目标策略网络 $\theta_e^- \leftarrow \tau \theta_o + (1 - \tau) \theta_e^-$
- 22: 软更新目标价值网络 $\omega_{i,e}^- \leftarrow \tau \omega_{i,o} + (1 - \tau) \omega_{i,e}^-$
- 23: **end if**
- 24: **end for**
- 25: **end for**

2.2.3 智能控制器部署

当 Actor 和 Critic 网络收敛, 则智能体训练完成。如图 2 所示, 训练好的 Actor 网络将作为神经网络控制策略在线实施, 并以端到端的方式生成高速飞行器

三通道姿态控制指令。具体来讲, 训练好的 Actor 网络接收到归一化后的姿态角跟踪误差和姿态角速度跟踪误差, 经过全连接层后, 输出三轴舵偏角。图 2 中的 x_{cmd} 为系统姿态角制导指令, x_g 为实际姿态角状态量, 即攻角 α 、侧滑角 β 、速度倾侧角 γ_v 。

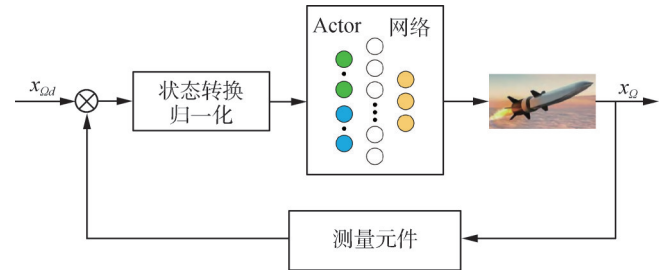


图2 神经网络模型的在线部署

Fig.2 Online deployment of neural network model

3 仿真试验及其结果分析

3.1 试验环境设置

通过仿真软件 Visual Studio Code, 编程语言 Python 对提出的 KG-TD3 模型进行训练和测试。Visual Studio Code 版本为 1.102.3, Python 版本为 3.12.3。用于试验的硬件平台配置如下: 操作系统为 Win11, CPU 为 Intel Core i5, 内存为 16 GB。

本研究选取高速飞行器的再入段, 其飞行持续时间为 38.4 s, 初始高度为 30 km, 初始速度约为 $Ma=5$, 动压范围从 27.4 kPa 到 588.8 kPa。图 3 为标准轨迹的典型参数, 从图 3 中可以看出, 该飞行阶段的高度、速度和动压呈现出显著而快速的变化, 同时伴随着模型参数的变化。飞行器机体参数如表 1 所示, 气动参数为参考文献 [17] 中提供的公开数据。

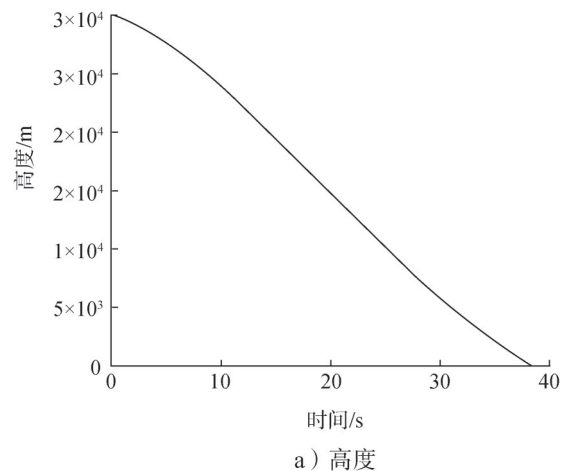
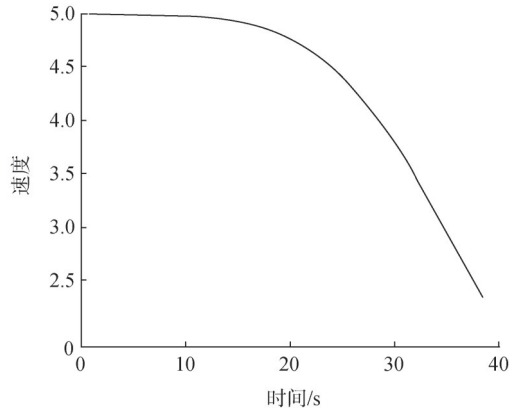
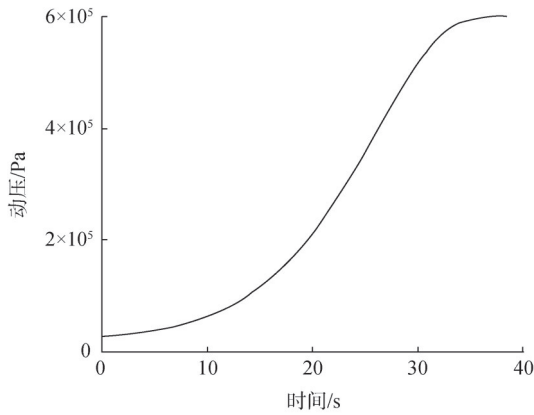


图3 标准轨迹的典型参数

Fig.3 Typical parameters of standard trajectories



b) 速度



c) 动压

续图3

表1 飞行器机体参数

Tab.1 Aircraft body parameters

参数	数值
m/kg	1 000
l/m	0.7
$J_y/(\text{kg}\cdot\text{m}^2)$	3 000
S/m^2	0.45
$J_x/(\text{kg}\cdot\text{m}^2)$	200
$J_z/(\text{kg}\cdot\text{m}^2)$	2 800

使用KG-TD3作为高速飞行器再入段三通道智能姿态控制器，其网络结构如表2所示，智能体训练过程中的超参数如表3所示。

表2 KG-TD3算法的网络结构

Tab.2 Network structure of the KG-TD3 algorithm

网络名称	层类型	神经元个数	激活函数
Actor	输入层	6	None
	全连接层	128	Relu
	输出层	3	Tanh
Critic	输入层	9	None
	全连接层	512	Relu
	全连接层	256	Relu
	输出层	1	Linear

表3 KG-TD3算法训练过程超参数

Tab.3 KG-TD3 algorithm training process hyper-parameters

超参数	数值
探索噪声 ϵ	$\mathcal{N}(0, 0.1^2)$
目标策略平滑噪声 ξ	$\mathcal{CN}(0, 0.2^2, -0.5, -0.5)$
延迟更新策略 k	2
目标网络更新率 τ	0.005
批量大小	128
Actor 学习率	$1 \times e^{-4}$
Critic 学习率	$1 \times e^{-4}$
行为克隆权重 λ	0.1

3.2 适应能力分析

为了验证本文提出的KG-TD3姿态控制器的适应能力，我们考虑了三个通道中的初始姿态偏差 $\Delta\alpha_0$ 、 $\Delta\beta_0$ 、 $\Delta\gamma_{v0}$ ；质量、气动参考面积、参考长度和大气密度偏差 Δm 、 ΔS 、 Δl 、 $\Delta\rho$ ；三通道中的惯性力矩偏差 ΔJ_x 、 ΔJ_y 、 ΔJ_z ；气动力系数偏差 ΔC_L 、 ΔC_Z ；俯仰力矩系数偏差 ΔC_{mz}^α 、 ΔC_{mz}^β 。表4为具体偏差值。由于风干扰不改变系统本身动力学模型结构，可通过添加扰动观测器等行为进行干扰估计并补偿。本研究重点针对需依赖控制器本身适应性的参数摄动，故适应能力分析只涉及参数摄动。

表4 参数偏差范围

Tab.4 Margins of error for parameters

参数	偏差范围
$\Delta\alpha_0/(\text{°})$	1
$\Delta\gamma_{v0}/(\text{°})$	3
ΔS	10%
$\Delta\rho$	20%
ΔJ_y	15%
ΔC_L	20%
ΔC_{mz}^α	20%
$\Delta\beta_0/(\text{°})$	1
Δm	10%
Δl	10%
ΔJ_x	15%
ΔJ_z	15%
ΔC_Z	20%
ΔC_{mz}^β	20%

图4展示了所提KG-TD3算法和TD3算法在训练过程中每回合智能体的动作奖励回报，反映了Actor网络的学习情况。KG-TD3算法在训练10轮左右奖励值趋近稳定，且无高位振荡的状态，反映了其训练过程稳定，收敛快速。相比之下，TD3算法在训练过程中奖励值波动较大，且数值略低于KG-TD3算法。

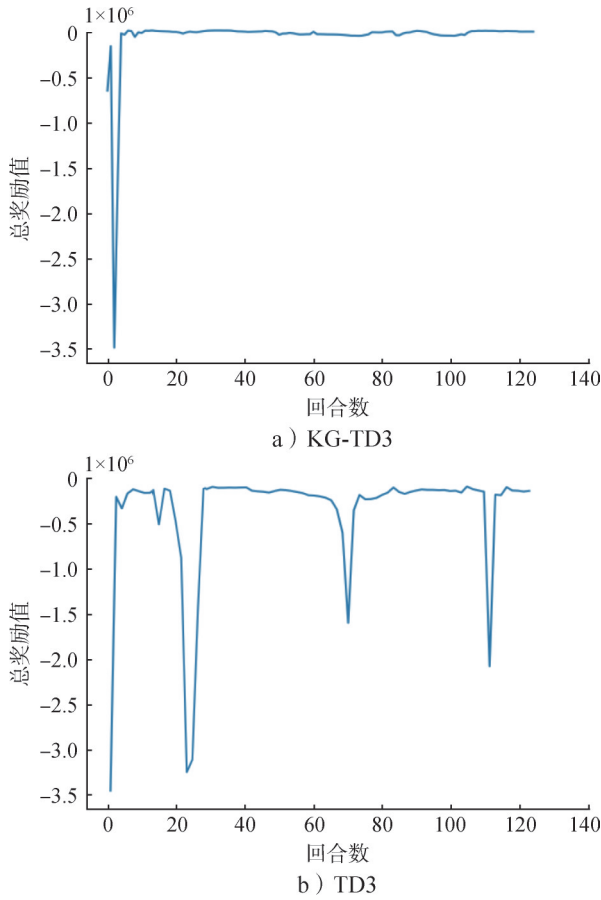


图4 训练过程中的奖励值

Fig.4 Reward progression during training

与TD3算法相比,本文提出的KG-TD3算法主要改进了奖励机制以及Actor网络优化更新策略。为了评估这一改进对姿态角跟踪精度的影响,我们将其与TD3进行对比分析,两种智能控制器的姿态角跟踪最大误差(MAX)如表5所示。其中,用于分析误差的飞行段为1~38.4 s。通过表5分析可知,KG-TD3相比TD3算法跟踪精度有较为明显的提升。

表5 姿态角跟踪误差对比

Tab.5 Comparison of attitude angle tracking errors

最大误差	TD3	KG-TD3
MAX _{α} (°)	2.403	0.420
MAX _{β} (°)	0.836	0.528
MAX _{γ} (°)	34.563	0.725

为验证所提控制方法的自适应能力,在三种典型工况下进行了仿真试验:标称状态、偏差上限状态和偏差下限状态。KG-TD3智能控制器的姿态角跟踪性能以及根据控制器得出的等效舵偏曲线如图5至图7所示。仿真试验结果分析表明,在14组涵盖上、下边界

的参数扰动工况下,系统稳态控制偏差均小于1°。

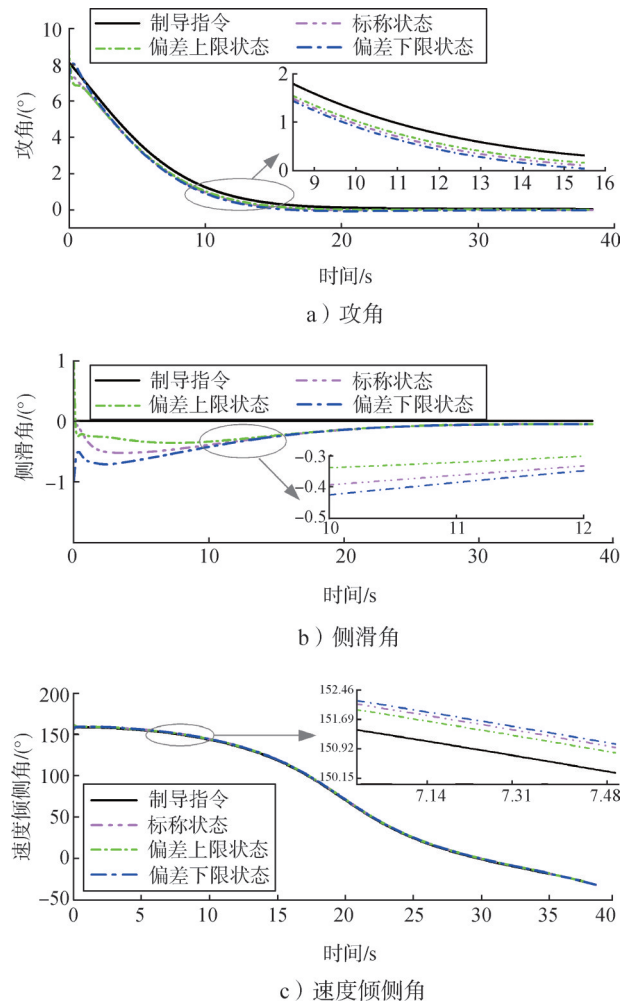


图5 姿态角跟踪效果

Fig.5 Attitude angle tracking effect

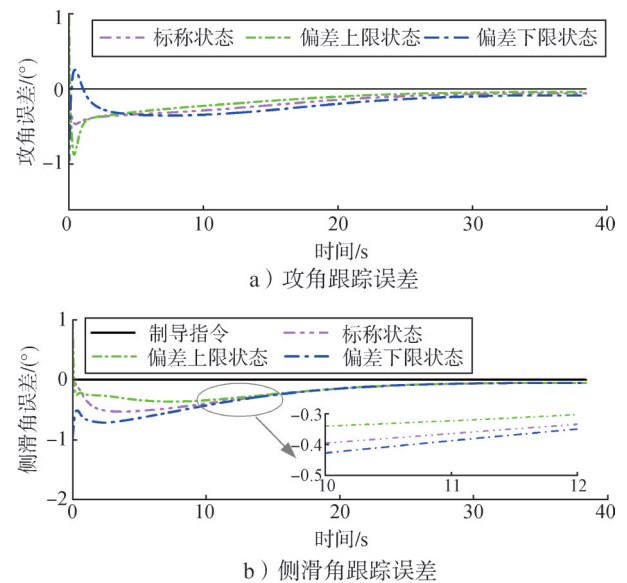


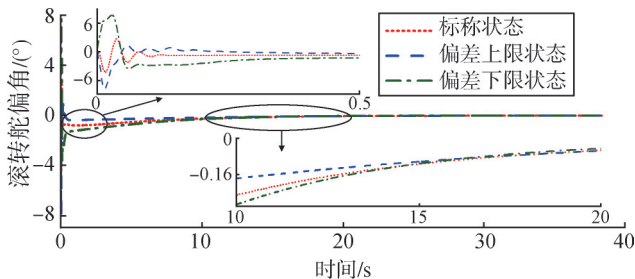
图6 姿态角跟踪误差

Fig.6 Attitude angle tracking error

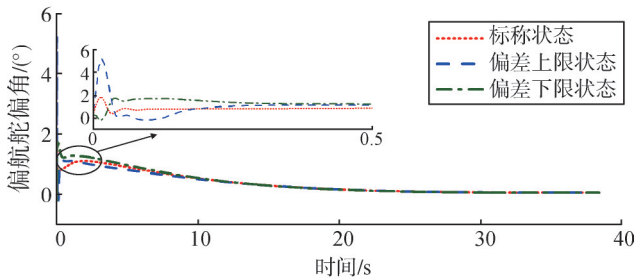


c) 速度倾侧角跟踪误差

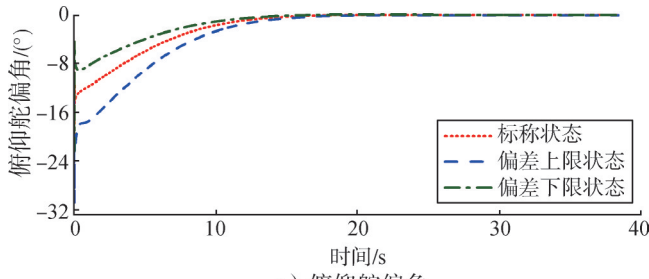
续图6



a) 滚转舵偏角



b) 俯仰舵偏角



c) 偏航舵偏角

图7 基于KG-TD3方法的三轴舵偏

Fig.7 Three-axis rudder deflection of the KG-TD3 method

参考文献

- [1] 刘双喜, 刘世俊, 李勇, 等. 国外高超声速飞行器及防御体系发展现状[J]. 空天防御, 2023, 6(3): 39-51.
LIU Shuangxi, LIU Shijun, LI Yong, et al. Current developments in foreign hypersonic vehicles and defense systems[J]. Air and Space Defense, 2023, 6(3): 39-51.
- [2] 樊铁, 秦昌茂, 董添, 等. 基于MIMO-ESO的高速飞行器自抗扰控制[J]. 导弹与航天运载技术(中英文), 2024(1): 64-70.
FAN Yi, QIN Changmao, DONG Tian, et al. ADRC attitude controller design for hypersonic vehicle based on MIMO-ESO[J]. Missiles and Space Vehicles, 2024(1): 64-70.
- [3] 包为民. 航天智能控制技术让运载火箭“会学习”[J]. 航空学报, 2021, 42(11): 8-17.
BAO Weimin. Space intelligent control technology enables launch vehicle to “self-learning” [J]. Acta Aeronautica et Astronautica Sinica, 2021, 42(11): 8-17.
- [4] ZHANG Z Y, MO Z B, CHEN Y T, et al. Reinforcement learning behavioral control for nonlinear autonomous system[J]. IEEE-CAA Journal of Automatica Sinica, 2022, 9(9): 1561-1573.
- [5] LUO B, SUN J Y, TANG R, et al. Reinforcement learning-based 3D trajectory tracking control of hypersonic gliding vehicles with time-varying uncertainties[J]. IEEE Transactions on Automation Science and Engineering, 2025(22): 8187-8199.
- [6] LIU C, DONG C Y, ZHOU Z J, et al. Barrier Lyapunov function based reinforcement learning control for air-breathing hypersonic vehicle with variable geometry inlet[J]. Aerospace Science and Technology, 2020(96): 105537.
- [7] LU K F, WANG W L, LIU X D, et al. Research progress and prospect of high-speed vehicle control technology based on reinforcement learning[J]. Advances in Astronautics, 2025, 8(2): 201-209.
- [8] WANG G, AN H, WANG Y, et al. Intelligent control of air-breathing hypersonic vehicles subject to path and angle-of-attack constraints[J]. Acta Astronautica, 2022(198): 606-616.
- [9] GAO Q, LI X, JI Y, et al. Research on active disturbance rejection control of hypersonic vehicle based on Q-learning[J]. Control Engineering of China, 2024, 31(4): 577-582.
- [10] LI X, JI Y H, SONG Y, et al. Modified deep deterministic policy gradient based on active disturbance rejection control for hypersonic vehicles[J]. Neural Computing and Applications, 2024, 36(8): 4071-4081.
- [11] 路坤锋, 贾晨辉, 黄旭, 等. 面向变构型飞行器的强化学习位置姿态一体化控制方法[J]. 宇航学报, 2024, 45(7): 1100-1110.

4 结束语

本研究针对高速飞行器再入段强非线性、高不确定性和参数快时变等复杂控制问题, 提出了一种基于改进型TD3算法的端到端智能姿态控制方法。通过融合混合奖励机制和基于行为克隆的先验知识约束, 有效解决了传统深度强化学习在姿态控制中训练不稳定、收敛困难的问题。

- LU Kunfeng, JIA Chenhui, HUANG Xu, et al. Reinforcement learning-based integrated position and attitude control method towards morphing flight vehicles[J]. *Journal of Astronautics*, 2024, 45(7): 1100-1110.
- [12] 姜凌峰, 李新凯, 张海, 等. 基于改进 TD3 算法的无人机动态环境无地图导航[J]. *航空学报*, 2025, 46(8): 298-313.
- JIANG Lingfeng, LI Xinkai, ZHANG Hai, et al. Mapless navigation of UAVs in dynamic environments based on an improved TD3 algorithm[J]. *Acta aeronautica et Astronautica Sinica*, 2025, 46(8): 298-313.
- [13] 彭博, 王晓波, 魏祥麟, 等. 基于 SPER-TD3 的无人机编队三维航迹规划[J]. *计算机系统应用*, 2025, 34(2): 61-73.
- PENG Bo, WANG Xiaobo, WEI Xianglin, et al. 3D trajectory planning for unmanned aerial vehicle formation based on SPER-TD3[J]. *Computer Systems & Applications*, 2025, 34(2): 61-73.
- [14] 闫雷鸣, 刘健, 朱永昕. DPC-DQRL: 动态行为克隆约束的离线-在线双 Q 值强化学习[J]. *计算机应用研究*, 2025, 42(4): 1003-1010.
- YAN Leiming, LIU Jian, ZHU Yongxin. DPC-DQRL: offline to online double Q value reinforcement learning with dynamic behavior cloning constraints[J]. *Application Research of Computers*, 2025, 42(4): 1003-1010.
- [15] 刘晓东, 黄万伟, 禹春梅. 含扩张状态观测器的高超声速飞行器动态面姿态控制[J]. *宇航学报*, 2015, 36(8): 916-922.
- LIU Xiaodong, HUANG Wanwei, YU Chunmei. Dynamic surface attitude control for hypersonic vehicle containing extended state observer[J]. *Journal of Astronautics*, 2015, 36(8): 916-922.
- [16] 黄旭, 柳嘉润, 张远, 等. 知识与数据混合驱动的高速飞行控制方法综述[J]. *宇航学报*, 2023, 44(8): 1113-1126.
- HUANG Xu, LIU Jiarun, ZHANG Yuan, et al. Review on knowledge-based and data-driver cooperating control methods of high-speed vehicle[J]. *Journal of Astronautics*, 2023, 44(8): 1113-1126.
- [17] LIU X, HUANG W, DU L. An integrated guidance and control approach in three-dimensional space for hypersonic missile constrained by impact angles[J]. *ISA Transactions*, 2017(66): 164-175.

作者简介

王伟丽 (1997—), 女, 博士研究生, 主要研究方向为飞行器智能控制、自适应控制等。

黄万伟 (1970—), 男, 博士, 研究员, 主要研究方向为飞行器制导与控制、智能控制、自适应控制等。

刘晓东 (1987—), 男, 博士, 研究员, 主要研究方向为飞行器制导与控制、智能控制、自适应控制等。

路坤锋 (1983—), 男, 博士, 研究员, 主要研究方向为飞行器制导与控制、智能控制、自适应控制等。

贾晨辉 (1985—), 男, 博士, 高级工程师, 主要研究方向为飞行器制导与控制、智能控制等。