

# 输入受限下的超紧密航天器编队避撞相对位置强化学习控制

孟亦真<sup>1,2</sup>, 黄静<sup>1,2†</sup>, 周绍辉<sup>3</sup>, 周彬<sup>1,2</sup>, 朱康武<sup>1,2</sup>

(1. 上海航天控制技术研究所, 上海 201109; 2. 上海市空间智能控制技术重点实验室, 上海 201109;

3. 上海航天空间技术有限公司, 上海 201109)

**摘要:** 考虑具有外界干扰、避撞约束和固定时间约束的近地轨道超紧密航天器编队的重构控制问题, 本文提出一种多约束条件下的考虑执行机构死区效应的航天器编队鲁棒控制方法。首先, 建立近地轨道完整的编队航天器相对位置非线性动力学方程和执行机构死区动态响应模型; 其次, 根据状态约束条件设计编队相对位置约束机制, 基于反步法和强化学习评判-动作网络, 提出防避撞约束和固定时间约束的鲁棒控制律, 进一步考虑到执行机构电推力器的死区效应, 基于强化学习的动作网络来近似死区特性, 本文通过最小化评判网络代价函数来解决执行机构死区效应对控制精度带来的影响, 并应用Lyapunov稳定性定理证明其闭环系统的一致有界性; 最后, 在MATLAB/Simulink平台上进行仿真验证, 结果表明所提出方法的有效性。

**关键词:** 航天器编队; 避撞; 强化学习; 死区效应; 固定时间约束

**引用格式:** 孟亦真, 黄静, 周绍辉, 等. 输入受限下的超紧密航天器编队避撞相对位置强化学习控制. 控制理论与应用, 2025, 42(4): 659 – 668

DOI: 10.7641/CTA.2023.20792

## Reinforcement learning control of collision avoidance for ultra-close formation of spacecraft with input constraints

MENG Yi-zhen<sup>1,2</sup>, HUANG Jing<sup>1,2†</sup>, ZHOU Shao-hui<sup>3</sup>, ZHOU Bin<sup>1,2</sup>, ZHU Kang-wu<sup>1,2</sup>

(1. Shanghai Aerospace Control Technology Institute, Shanghai 201109, China;

2. Shanghai Key Laboratory of Aerospace Intelligent Control Technology, Shanghai 201109, China;

3. Shanghai Space Technology Co., Ltd, Shanghai 201109, China)

**Abstract:** Considering the control problem of reconstructing the ultra-tight formation of near-Earth orbit spacecraft in the presence of external disturbances, collision avoidance constraints, and fixed-time constraints, this study presents a robust control method for spacecraft formation that accounts for the dead-zone effect of the actuator under multiple constraint conditions. Firstly, we establish the nonlinear dynamic equations governing the relative positions of the spacecraft in the complete near-Earth orbit formation, as well as the dynamic response model capturing the dead-zone effect of the actuator. Secondly, we design a constraint mechanism for the relative positions of the formation based on state constraints. Robust control laws, employing a combination of backstepping and a reinforcement learning actor-critic network, are proposed to address collision avoidance constraints and fixed-time constraints. Additionally, we approximate the dead-zone characteristics of the actuator's thrusters by leveraging a reinforcement learning actor network. To mitigate the impact of the dead-zone effect on control accuracy, we minimize the cost function of the actor network. The Lyapunov stability theorem is employed to demonstrate the uniformly boundedness of the closed-loop system. Finally, we conduct simulation verification on the MATLAB/Simulink platform, and the results substantiate the effectiveness of the proposed method.

**Key words:** spacecraft formation; collision avoidance; reinforcement learning control; dead-zone effect; fixed time constraint

**Citation:** MENG Yizhen, HUANG Jing, ZHOU Shaohui, et al. Reinforcement learning control of collision avoidance for ultra-close formation of spacecraft with input constraints. *Control Theory & Applications*, 2025, 42(4): 659 – 668

收稿日期: 2022-09-08; 录用日期: 2023-12-07.

†通讯作者. E-mail: huangjing04415@163.com.

本文责任编辑: 贾英民.

国家重点研发计划项目(2022YFB3902700, 2022YFB3902702), 空间目标感知全国重点实验室资助.

Supported by the National Key R&D Program of China (2022YFB3902700, 2022YFB3902702) and the National Key Lab of Space Target Awareness.

## 1 引言

航天器编队飞行因其结构灵活、功能强大、可靠性高、生命周期长及发射风险低等优点,已成为航天控制领域的研究热点.与单个航天器任务相比,航天器编队飞行面临的挑战是在较长时间内维持固定的航天器相对距离,如深空干涉仪器、合成孔径雷达等.这些复杂、精密的编队任务需对编队中各个航天器到达编队构型的时间、相对位置、控制精度进行约束,以满足指定的一致性的任务需求(形成新构型、波束同步等).由此,便涉及到航天器系统状态约束、控制时间约束、分布式协同与电推进执行机构的智能控制问题<sup>[1]</sup>.

有别于单个航天器任务,以主-从结构编队组成的分布式雷达天线为例,为实现高分辨率的观测,天线阵面边缘之间的距离必须保持在1米之内.为实现这种超近距离、高精度的航天器编队飞行任务,如何快速形成编队构型,同时避免航天器间的碰撞,并能满足精度需求,含有死区约束电推进下的分布式状态约束、避碰控制系统的设计俨然成为超紧密航天器编队重构任务的重中之重.

在航天器超紧密编队任务中,其超近距离的重构,涉及到航天器的系统状态约束控制以及受到外界扰动后的避碰问题.针对于系统状态约束控制问题,已有众多的研究成果,如基于模型预测控制方法<sup>[2]</sup>的不变集方法<sup>[3]</sup>、障碍李雅普诺夫函数以及积分型障碍李雅普诺夫函数<sup>[4]</sup>.文献[5]采用障碍Lyapunov函数结合反步法解决航天器预置性能下的状态约束问题.而目前基于反步法结合障碍Lyapunov函数的方法在超紧密编队的全状态约束控制设计中,还需要满足虚拟控制律的可行性条件,即其必须要满足给定的约束区<sup>[6]</sup>.为克服该预先给定约束域带来的应用限制,本文提出一个仅依赖于系统状态的非线性约束函数,以此来约束紧密航天器编队的相对位置,实现控制目标.同时,因为超紧密编队,在系统受到扰动时,避免航天器间碰撞与系统约束状态的越界成为超紧密编队重构控制中的重要一环.解决这个问题的常见方法是构造李雅普诺夫势函数,采用李雅普诺夫函数完成协同控制和障碍势函数进行避碰,随后,设计由协调项和避碰项组成的集成控制器,完成控制目标<sup>[7]</sup>.该方法难以实现在避免碰撞的同时,保持系统的相对位置状态约束.为解决该问题,借鉴强化学习能够基于评判网络对当前的控制效果进行评价,而后通过动作网络的控制调节,使得系统的收益在未来最大化<sup>[6]</sup>.基于此,本文拟用强化学习的指标函数,来实现惩罚功能,当超紧密编队卫星进入碰撞或者发生系统状态越界行为后,强化学习评判此刻的控制效果差(输出较大的惩罚值),进而调节动作网络的输出,从而能够最小化

评判网络的输出值,迫使越界的航天器自动恢复到期望的约束区域,避免因碰撞或者越界造成的损坏及构型破坏.

另一个值得关注的问题是编队控制的时间约束问题,尤其是对于快速高精度控制系统,比如多卫星系统的控制、多导弹飞行控制等,要求各子系统在有限时间内能够达到一致性.在已有的有限时间控制问题中,收敛时间依赖于初值,从而无法解决初值未知或比较大时,难以满足系统编队的有限时间控制需求.为克服该不足,本文设计基于固定时间的控制策略,其收敛时间有不依赖于初值的上界,而收敛过程中的通讯距离由航天器的状态约束来保持<sup>[8]</sup>.同时,为保障航天器编队相对位置控制系统的精度,电推进因其比冲较化学推进至少高一个量级,同时兼具精度可调、冲量小且寿命长的优点,成为航天器编队相对位置控制适宜的执行机构,其能够提升航天器任务执行能力,扩展其空间任务范围<sup>[9]</sup>.但电推进器的死区效应给控制系统设计带来了负面影响.

针对多约束条件下的超紧密编队的控制问题,本文的创新点如下:

1) 基于固定时间收敛的非线性控制机制.在反步法的基础上,通过引入固定时间收敛的非线性控制项,解决超紧密编队分布式控制的有限时间约束问题.

2) 基于状态依赖的约束机制.考虑航天器超紧密编队相对位置的约束,提出基于状态依赖的约束机制,以确保编队的相对位置约束.

3) 基于强化学习的评判-动作网络的软约束控制策略.为解决航天器状态越界导致的碰撞和编队构型破坏问题,设计基于强化学习的评判-动作网络的软约束控制策略.通过最小化评判网络输出作为目标,并优化动作网络的输出,发挥其状态约束机制,确保航天器间的相对位置约束.

本文提出的基于强化学习的电推进鲁棒控制算法创新地解决多约束条件下的超紧密编队的控制问题.通过引入固定时间收敛的非线性控制项、基于状态依赖的约束机制和基于强化学习的评判-动作网络的软约束控制策略,有效地应对约束条件和执行机构死区效应带来的挑战.此外,通过稳定性分析和MATLAB/Simulink平台上的仿真验证,证明该方法的实时性和有效性.

## 2 编队航天器相对位置数学模型

本文采用主-从方式建立编队航天器相对位置数学模型,假设每个航天器都是一个质点,领航者在以恒定角速度运行在圆形轨道上,而跟随者在领航者的周围按照所期望的参考轨迹运行.为获得航天器编队飞行(spacecraft formation flying, SFF)的运动模型,结合参考文献[10],引入如图1所示的坐标系:

1)  $C_1 = \{X, Y, Z\}$  表示以地心为原点的惯性坐标系. 定义  $\mathbf{R}$  为该系下从地心指向领航者的位置矢量.  $\boldsymbol{\rho} \in \mathbb{R}^3$  为主航天器与从航天器间的相对位置矢量.

2)  $C_L = \{x_1, y_1, z_1\}$  为主航天器参考坐标系,  $x_1$  的正方向与瞬时的切向速度相反;  $Y_1$  的正方向为沿着矢径  $\mathbf{R}$  的方向;  $Y_1$  垂直于  $X_1$  和  $Z_1$  组成的平面,  $x_1, y_1, z_1$  符合右手准则.

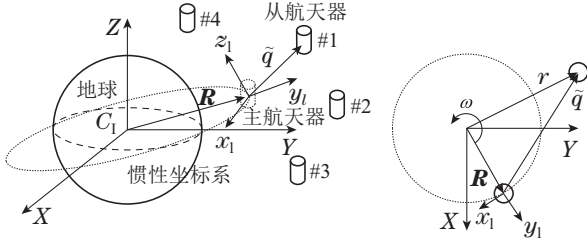


图1 航天器编队运行示意图

Fig. 1 Schematic representation of the SFF system

在  $C_1$  坐标系下, 主航天器和从航天器之间相对位置运动的动力学方程如下<sup>[1]</sup>:

$$m_f \ddot{\boldsymbol{\rho}} + m_f M_G G \left( \frac{\mathbf{R} + \boldsymbol{\rho}}{\|\mathbf{R} + \boldsymbol{\rho}\|^3} - \frac{\mathbf{R}}{\|\mathbf{R}\|^3} \right) + \mathbf{F}_d = \mathbf{u}_t, \quad (1)$$

其中:  $m_1$  和  $m_f$  分别表示主航天器和从航天器的质量;  $G$  为万有引力常数;  $M_G$  表示地球的质量;  $\mathbf{F}_d \triangleq \mathbf{F}_{df} - \left(\frac{m_f}{m_1}\right) \mathbf{F}_{dl}$  为编队系统受到的扰动量,  $\mathbf{F}_{df}, \mathbf{F}_{dl}$  分别表示作用在主航天器和从航天器上的外力干扰;  $\mathbf{u}_t = \mathbf{u}_f - \left(\frac{m_f}{m_1}\right) \mathbf{u}_l$  为系统的控制量,  $\mathbf{u}_f, \mathbf{u}_l$  分别表示作用在主航天器和从航天器上的控制力. 进一步, 将位置矢量  $\boldsymbol{\rho}$  映射到动坐标系  $C_L$  中, 可得

$$\boldsymbol{\rho} = x_1 \vec{i}_1 + y_1 \vec{j}_1 + z_1 \vec{k}_1, \quad (2)$$

其中  $\vec{i}_1, \vec{j}_1, \vec{k}_1$  表示  $C_L$  三坐标轴的单位矢量, 那么它的二阶导数方程可以改写为

$$\ddot{\boldsymbol{\rho}} = (\ddot{x} - 2\omega_0 \dot{y} - \omega_0^2 x) \vec{i}_1 + (\ddot{y} + 2\omega_0 \dot{x} - \omega_0^2 y) \vec{j}_1 + \ddot{z} \vec{k}_1, \quad (3)$$

进一步, 将式(3)代入式(1)中, 可得跟随者相对于领航者的非线性相对位置平动动力学方程为

$$m_f \ddot{\boldsymbol{\rho}} + m_f \mathbf{C}(\omega) \dot{\boldsymbol{\rho}} + m_f \mathbf{N}(\boldsymbol{\rho}, \omega, \mathbf{R}, \mathbf{u}_1) = \mathbf{u}_f, \quad (4)$$

其中:  $\mathbf{C}(\omega_0) = 2\omega_0 \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{3 \times 3}$ ,  $\mathbf{N}(\cdot)$  为

$$\mathbf{N}(\boldsymbol{\rho}, \omega, \mathbf{R}, \mathbf{u}_1) = M_G G \left[ \frac{x}{\|\mathbf{R} + \boldsymbol{\rho}\|^3} - \frac{\omega_0^2 y z}{\|\mathbf{R} + \boldsymbol{\rho}\|^3} - \omega_0^2 x \left( \frac{y + \|\mathbf{R}\|}{\|\mathbf{R} + \boldsymbol{\rho}\|^3} - \frac{1}{\|\mathbf{R}\|^2} \right) \right]^T.$$

为了便于后续的推导, 基于式(1)(4), 可得SFF的模型为

$$\begin{cases} \dot{\mathbf{x}}_1 = \mathbf{f}_1(\bar{\mathbf{x}}_1, \mathbf{p}_1) + \mathbf{g}_1(\bar{\mathbf{x}}_1) \mathbf{x}_2, \\ \dot{\mathbf{x}}_2 = \mathbf{f}_2(\bar{\mathbf{x}}_2, \mathbf{p}_2) + \mathbf{g}_2(\bar{\mathbf{x}}_2) \mathbf{N}(\mathbf{v}(t)) + \mathbf{D}, \end{cases} \quad (5)$$

其中:  $\mathbf{x}_1 = \boldsymbol{\rho}, \mathbf{x}_2 = \dot{\boldsymbol{\rho}}, \mathbf{f}_1(\bar{\mathbf{x}}_1, \mathbf{p}_1) = 0, \mathbf{g}_1(\bar{\mathbf{x}}_1) = \mathbf{I}_{3 \times 3}, \mathbf{f}_2(\bar{\mathbf{x}}_2, \mathbf{p}_2) = -\mathbf{M}_f^{-1} \mathbf{C}(\mathbf{x}_2) \mathbf{x}_2 - \mathbf{M}_f^{-1} \mathbf{N}(\mathbf{x}_1), \mathbf{g}_2(\bar{\mathbf{x}}_2) = \mathbf{M}_f^{-1}$ .

航天器受到的干扰  $\mathbf{D} = -\mathbf{M}_f^{-1} \mathbf{F}_d$ .  $\mathbf{N}(\mathbf{v}(t))$  为系统的控制输入量. 由于受限于与小卫星的重量、体积等因素, 具有重量小、比冲高、电推进器成为各国争相研究的重点. 在本文中, 采用电推力器作为航天器编队控制的执行机构, 重点考虑执行机构的死区特性, 定义为系统带有死区约束的输入控制量, 为执行机构提供的控制量,  $\mathbf{N}(\mathbf{v}(t)) = [N_1(v_1(t)) \ N_2(v_2(t)) \ N_3(v_3(t))]^T$  具体表达如下<sup>[11-12]</sup>:

$$N_i(v_i(t)) = \begin{cases} h_r(v_i(t) - o_r), & v_i(t) > o_r, \\ 0, & |v_i(t)| \leq o_r, \\ h_l(v_i(t) + o_l), & v_i(t) < -o_l, \end{cases} \quad (6)$$

其中:  $i = 1, 2, 3$ ;  $h_l, h_r$  分别为死区的左、右斜率;  $o_l, o_r$  分别为执行器的非线性断点. 经分析,  $\mathbf{N}(\mathbf{v}(t)) = \mathbf{v}(t) - \Delta \mathbf{v}(t)$ ,  $\Delta \mathbf{v}(t)$  为死区效应带来的执行误差.

### 3 固定时间强化学习优化控制

#### 3.1 航天器状态约束

考虑到编队航天器之间的避撞需求, 本文将引入系统的非对称状态约束, 设计如下的非线性状态依赖函数:

$$\zeta_i(t) = \frac{x_{1,i}(t)}{(F_{i1}(t) + x_{1,i}(t))(F_{i2}(t) - x_{1,i}(t))}, \quad (7)$$

其中:  $i = 1, 2, 3$ ;  $\boldsymbol{\zeta}(t) = [\zeta_1(t) \ \zeta_2(t) \ \zeta_3(t)]^T$ ,  $\mathbf{x}_1(t) = [x_{11}(t) \ x_{12}(t) \ x_{13}(t)]^T$ ;  $F_{i1}(t)$  和  $F_{i2}(t)$  为系统状态约束函数, 并且系统初始值  $x_{1i}(0)$  满足  $-F_{i1}(0) < x_{1i}(0) < F_{i2}(0)$ . 同时, 式(7)的引入可以直接处理系统的全状态约束, 并且避免了对虚拟控制律可行性条件的判别. 为便于后续的推导分析, 本文基于反步法引入如下的坐标变换:

$$\begin{cases} \mathbf{z}_1 = \boldsymbol{\zeta}_1 - \boldsymbol{\alpha}_0, \mathbf{z}_2 = \mathbf{x}_2 - \boldsymbol{\alpha}_{2f}, \\ \boldsymbol{\alpha}_{0,i} = \frac{y_{d,i}}{(F_{i,1}(t) + y_{d,i})(F_{i,2}(t) - y_{d,i})}, \end{cases} \quad (8)$$

其中:  $i = 1, 2, 3$ ; 状态约束下的跟踪误差  $\mathbf{z}_1$  定义为  $\mathbf{z}_1 = [z_{1,1} \ z_{1,2} \ z_{1,3}]^T$ ;  $\boldsymbol{\alpha}_0 = [\alpha_{0,1} \ \alpha_{0,2} \ \alpha_{0,3}]^T$  为虚拟控制量;  $\boldsymbol{\alpha}_{2f} = [\alpha_{2f,1} \ \alpha_{2f,2} \ \alpha_{2f,3}]^T$  为避免维数爆炸, 引入如下的一阶滤波器状态:

$$\varepsilon_1 \dot{\boldsymbol{\alpha}}_{2f,i} + \boldsymbol{\alpha}_{2f,i} = \boldsymbol{\alpha}_{1,i}, \quad (9)$$

其中:  $i = 1, 2, 3$ ;  $\varepsilon_1 > 0$ ; 虚拟控制量  $\boldsymbol{\alpha}_1 = [\alpha_{1,1} \ \alpha_{1,2} \ \alpha_{1,3}]^T$  作为一阶惯性环节的输入. 进一步, 定义如下因一阶惯性环节引入引起的误差量  $\mathbf{y}_1$ , 具体为

$$\mathbf{y}_1 = \boldsymbol{\alpha}_{2f} - \boldsymbol{\alpha}_1, \quad (10)$$

其中  $\mathbf{y}_1 = [y_{11} \ y_{12} \ y_{13}]^T$ .

下面基于反步法思想,分步骤详细阐述控制器算法的设计过程:

**步骤 1** 对  $\zeta_1$  关于时间求导,可得

$$\dot{\zeta}_1 = \boldsymbol{\eta}_1 \dot{\mathbf{x}}_1, \quad (11)$$

因此,结合式(11),可得  $\mathbf{z}_1 = \zeta_1 - \boldsymbol{\alpha}_0$  关于时间  $t$  的导数

$$\dot{\mathbf{z}}_1 = \boldsymbol{\eta}_1 (\mathbf{g}_1(\mathbf{x}_1) \mathbf{x}_2 + \mathbf{f}_1) - \boldsymbol{\eta}_{d1} \dot{\mathbf{y}}_d, \quad (12)$$

其中:

$$\boldsymbol{\eta}_{d1} = [\eta_{d1,1} \ \eta_{d1,2} \ \eta_{d1,3}]^T, \\ \eta_{d1,i} = \frac{F_{i,1} F_{i,2} + y_{d,i}^2}{(F_{i,1} + y_{d,i})^2 (F_{i,2} - y_{d,i})^2}, \quad i = 1, 2, 3.$$

结合  $\mathbf{z}_2$  与  $\mathbf{y}_1$  的定义,可得  $\mathbf{x}_2 = \mathbf{z}_2 + \mathbf{y}_1 + \boldsymbol{\alpha}_1$ . 进一步,对  $\frac{1}{2} \mathbf{z}_1^T \mathbf{z}_1$  求关于时间的  $t$  的导数,可得

$$\mathbf{z}_1^T \dot{\mathbf{z}}_1 = \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{g}_1 \boldsymbol{\alpha}_1 + \mathcal{P}_1, \quad (13)$$

其中  $\mathcal{P}_1 = \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{f}_1 + \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{g}_1 (\mathbf{z}_2 + \mathbf{y}_1) - \mathbf{z}_1^T \boldsymbol{\eta}_{d1} \dot{\mathbf{y}}_d$ . 根据如下的杨氏不等式:

$$\mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{f}_1 \leq \underline{g}_1 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{f}_1\|^2 + \frac{1}{4\underline{g}_1}, \quad \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{g}_1 \mathbf{z}_2 \leq \underline{g}_2 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_2\|^2 + \frac{\bar{g}_1^2}{4\underline{g}_2}, \quad \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{g}_1 \mathbf{y}_1 \leq \underline{g}_1 \frac{\bar{g}_1^2}{\underline{g}_1} \|\mathbf{z}_1\|^2 \times \|\boldsymbol{\eta}_1\|^2 + \frac{1}{4} \|\mathbf{y}_1\|^2, \quad -\mathbf{z}_1^T \boldsymbol{\eta}_{d1} \dot{\mathbf{y}}_d \leq \underline{g}_1 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_{d1}\|^2 \|\dot{\mathbf{y}}_d\|^2 + \frac{1}{4\underline{g}_1},$$
 因此,  $\mathcal{P}_1$  可化简为

$$\mathcal{P}_1 \leq \underline{g}_1 b_1 \|\mathbf{z}_1\|^2 \Pi_1 + \Delta_1 + \underline{g}_2 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_2\|^2 + \frac{\|\mathbf{y}_1\|^2}{4}, \quad (14)$$

其中:  $b_1 = \max\{1, \|\mathbf{f}_1\|^2, \frac{\bar{g}_1^2}{\underline{g}_1}\}$ ;  $\underline{g}_1, \underline{g}_2$  为正定参数;  $\Pi_1 = 2\|\boldsymbol{\eta}_1\|^2 + \|\boldsymbol{\eta}_{d1}\|^2 \|\dot{\mathbf{y}}_d\|^2$  已知,同时为可计算的标量函数;  $\Delta_1 = \frac{1}{2\underline{g}_1} + \frac{\bar{g}_1^2}{4\underline{g}_2}$  为未知的正常数;  $\underline{g}_1, \underline{g}_2 > 0$  代表控制分配矩阵范数的下界,同理  $\bar{g}_1, \bar{g}_2 > 0$  指代的是控制分配矩阵的上界.

进一步,可将式(13)化简为

$$\mathbf{z}_1^T \dot{\mathbf{z}}_1 \leq \mathbf{z}_1^T \boldsymbol{\eta}_1 \mathbf{g}_1 \boldsymbol{\alpha}_1 + \underline{g}_1 b_1 \|\mathbf{z}_1\|^2 \Pi_1 + \Delta_1 + \underline{g}_2 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_2\|^2 + \frac{1}{4} \|\mathbf{y}_1\|^2, \quad (15)$$

从而,设计如下的虚拟控制律  $\boldsymbol{\alpha}_1$ :

$$\begin{cases} \boldsymbol{\alpha}_1 = -\text{diag}^{-1}\{\boldsymbol{\eta}_1\} (\mathbf{k}_1 \mathbf{z}_1 + \mathbf{z}_1 \hat{b}_1 \Pi_1 + l_1 \text{sig}^{\frac{m_1}{n_1}}(\mathbf{z}_1) + l_2 \text{sig}^{\frac{p_1}{r_1}}(\mathbf{z}_1)), \\ \dot{\hat{b}}_1 = \ell_1 (\|\mathbf{z}_1\|^2 \Pi_1 - \ell_{b_1} \hat{b}_1), \end{cases} \quad (16)$$

其中:  $\mathbf{k}_1 = \text{diag}\{k_{11}, k_{22}, k_{33}\}$ ,  $k_{11}, k_{22}, k_{33}, \ell_1, \ell_{b_1}$  为

待设计的正常数;  $\tilde{b}_1 = b_1 - \hat{b}_1$  为估计误差,  $\hat{b}_1$  为  $b_1$  的估计值. 函数  $\text{sig}^{\frac{p_1}{r_1}}(\cdot) = \text{sgn}(\cdot) |\cdot|^{\frac{p_1}{r_1}}$ .

在进行对虚拟控制量  $\boldsymbol{\alpha}_1$  的稳定性分析之前,结合  $\mathbf{y}_1$  的定义(10),可得  $\dot{\mathbf{y}}_1 = -\frac{\mathbf{y}_1}{\varepsilon_1} + \dot{\mathbf{h}}_1$ ,  $\dot{\mathbf{h}}_1 = -\dot{\boldsymbol{\alpha}}_1 = -\frac{\partial \boldsymbol{\alpha}_1}{\partial \boldsymbol{\eta}_1} \dot{\boldsymbol{\eta}}_1 - \frac{\partial \boldsymbol{\alpha}_1}{\partial \mathbf{z}_1} \dot{\mathbf{z}}_1 - \frac{\partial \boldsymbol{\alpha}_1}{\partial \hat{b}_1} \dot{\hat{b}}_1 - \frac{\partial \boldsymbol{\alpha}_1}{\partial \Pi_1} \dot{\Pi}_1$ . 因此可得  $\mathbf{y}_1^T \dot{\mathbf{y}}_1 = -\frac{\mathbf{y}_1^T \mathbf{y}_1}{\varepsilon_1} + \mathbf{y}_1^T \dot{\mathbf{h}}_1 \leq (\frac{1}{4} - \frac{1}{\varepsilon_1}) \mathbf{y}_1^T \mathbf{y}_1 + \dot{\mathbf{h}}_1^T \mathbf{h}_1$ .

而后,通过选取如下的Lyapunov函数:  $V_1 = \frac{1}{2} \mathbf{z}_1^T \times \mathbf{z}_1 + \frac{g_1}{2\ell_{b_1}} \tilde{b}_1^2 + \frac{1}{2} \mathbf{y}_1^T \mathbf{y}_1$ , 并对其求关于时间  $t$  的导数

$$\dot{V}_1 \leq -\underline{g}_1 \lambda_{\min}(\mathbf{k}_1) \|\mathbf{z}_1\|^2 - \frac{g_1 \ell_{b_1}}{2} \|\tilde{b}_1\|^2 - \bar{\varepsilon}_1 \mathbf{y}_1^T \mathbf{y}_1 + \dot{\mathbf{h}}_1^T \mathbf{h}_1 + \underline{g}_2 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_2\|^2 + \Gamma_1 - \mathbf{z}_1^T \underline{g}_1 (l_1 \text{sig}^{\frac{m_1}{n_1}}(\mathbf{z}_1) + l_2 \text{sig}^{\frac{p_1}{r_1}}(\mathbf{z}_1)), \quad (17)$$

其中:  $0 < \underline{g}_1 \leq g_1$ ,  $\underline{g}_1 \ell_{b_1} \tilde{b}_1^2 \leq -\frac{g_1 \ell_{b_1}}{2} \|\tilde{b}_1\|^2 + \frac{g_1 \ell_{b_1}}{2} \times \|\tilde{b}_1\|^2$ ,  $\Gamma_1 = \Delta_1 + \frac{g_1 \ell_{b_1}}{2} \|\tilde{b}_1\|^2$ .  $\underline{g}_2 \|\mathbf{z}_1\|^2 \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_2\|^2$  将会在后期的控制器设计中消除,此外,  $\dot{\mathbf{h}}_1$  将会在后期的控制算法的最后一项中处理.

### 3.2 强化学习评判网络设计

结合控制框图(如图2所示),为满足系统的状态约束,提高系统的安全性,本节在第3.1节状态约束的基础上,设计如下的效用函数为软约束:

$$\mathbf{J}(t) = \int_t^\infty \gamma^{-\frac{\xi-t}{T}} \mathbf{p}(\mathbf{s}(\xi)) d\xi, \quad (18)$$

其中:  $\xi$  为关于时间变量  $t$  的积分变量;  $T > 0$  为强化学习效用函数的解算区间;  $\gamma \in (0, 1)$  为折扣因子,用来减小当前权重对未来系统收益的影响,更符合客观事实规律. 如果  $\mathbf{s}$  足够小,则达到系统的优化目标.

1) 将受到干扰的系统状态重新约束到硬约束范围内;

2) 优化电推进器的控制输出,使得系统的能量消耗越小越好,即强化学习的效用函数不会再持续增长,如果  $\mathbf{s}$  比较大,则自适应调节控制器输出使得状态约束下的跟踪误差  $\mathbf{z}_1$ ,  $\dot{\mathbf{z}}_1$  和  $\mathbf{J}(t)$  减小.

本文的控制目标主要是使状态约束下的跟踪误差  $\mathbf{z}_1$ , 强化学习的效用函数  $\mathbf{p}(\mathbf{s})$ , 状态量  $\mathbf{s}$  足够小的同时保障其有界. 进而,容易获得强化学习指标函数的期望值  $\mathbf{J}_d = [0 \ 0 \ 0]^T$ . 根据参考文献[13]的设计如下:

$$p_i(\mathbf{s}_i(\xi)) = \begin{cases} 0, & F_{i1}(t) < x_{1i} < F_{i2}(t), \quad s_i(v(t)) \leq \varepsilon \\ \varepsilon, & F_{i1}(t) < x_{1i} < F_{i2}(t), \quad s_i(v(t)) > \varepsilon, \\ \varphi, & x_{1i} \geq F_{i2}(t), \quad x_{1i} \leq F_{i1}(t), \end{cases} \quad (19)$$

其中  $\mathbf{s}(v(t)) = \mathbf{z}_1^T \mathbf{R}_s \mathbf{z}_1 + \mathbf{v}^T \mathbf{Q}_s \mathbf{v}$ . 为实现上述两项

优化目标, 强化学习目标函数 $p_i(s_i(\xi))$ 与优化目标的关系分析如下:

1) 当系统相对距离约束在界限 $F_{i1}$ 内, 同时, 系统的能量消耗小于设定阈值 $\varepsilon$ 时, 判定系统的控制效果较好, 使得 $p_i(s_i(\xi)) = 0$ , 表明此时有效补偿电推进器死区特性的影响;

2) 当 $p_i(s_i(\xi)) = \varepsilon$ 时, 表明此时的死区特性控制补偿量并未达到较好的效果, 需通过调节控制输出, 减小效用函数的输出, 以此来优化死区补偿控制量的输出. 在最小化效用函数的过程中, 有效解决电推进器的死区效应, 完成目标2;

3) 当 $p_i(s_i(\xi)) = \varphi$ 时,  $\varphi > 0$ 为一个足够大的常数, 当系统状态溢出约束边界后, 起到惩罚作用, 使得系统状态重新返回约束边界内, 实现优化目标1.

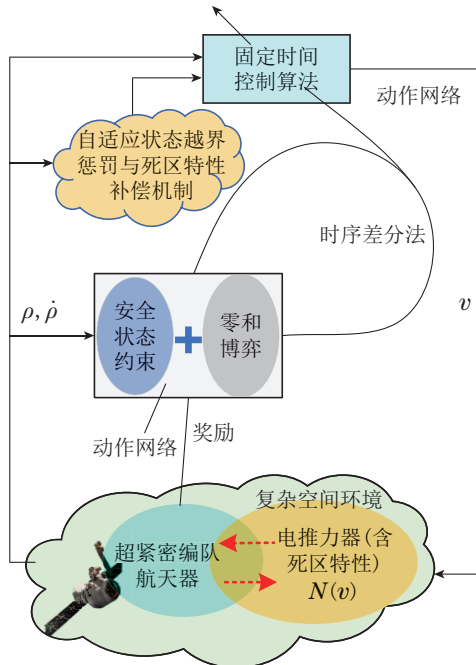


图2 基于强化学习的超紧密航天器编队防撞相对位置控制算法框图

Fig. 2 Diagram of relative position control for ultra-tight spacecraft formation collision avoidance based on reinforcement learning

依据参考文献[13], 强化学习的目标是使效用函数 $J(t)$ 越小越好. 结合式(18), 由于 $J(t)$ 的积分是到达无穷的, 为解决该优化指标无穷时域内的求解问题, 本文建如下的Bellman误差方程, 反映 $J(t-T)$ 与 $J(t)$ 的递推迭代关系:

$$J(t-T) = \int_{t-T}^{\infty} \gamma^{\frac{-\xi+t-T}{T}} \mathbf{p}(s(\xi)) d\xi = \gamma^{-1}(\mathbf{J}(t) + \mathbf{p}_c), \quad (20)$$

其中 $\mathbf{p}_c = \int_{t-T}^t \gamma^{\frac{(-\xi+1)}{T}} \mathbf{p}(s(\xi)) d\xi$ 为价值函数在积分区间 $[t-T, t)$ 的约束值.

进一步,  $\mathbf{p}_c = [p_{c1} \ p_{c2} \ p_{c3}]^T$ 可得

$$p_{ci} = \int_{t-T}^t \gamma^{\frac{-\xi+1}{T}} p_i(s_i(\xi)) d\xi = \begin{cases} 0, & F_{i1}(t) < x_{1,i} < F_{i2}(t), \ s(\mathbf{v}(t)) \leq \varepsilon, \\ \frac{\varepsilon(\gamma-1)}{\ln \gamma} T, & F_{i1}(t) < x_{1,i} < F_{i2}(t), \ s(\mathbf{v}(t)) > \varepsilon, \\ \frac{\varphi(\gamma-1)}{\ln \gamma} T, & x_{1,i} \geq F_{i2}(t), \ x_{1,i} \leq F_{i1}(t), \end{cases} \quad (21)$$

其中:  $i = 1, 2, 3$ ;  $\|\mathbf{p}_c\| \leq b_{p_c}$ ,  $b_{p_c} > 0$ 为一正常数.

基于参考文献[7], 采用径向基神经网络来逼近强化学习评判指标函数 $J(t)$ , 具体为

$$J(t) = \mathbf{W}_c^{*T} \phi_c(\mathbf{x}_c(t)) + \varepsilon_c(\mathbf{x}_c(t)), \quad (22)$$

其中:  $\mathbf{W}_c^*$ 为理想的权重, 存在正常数 $b_{W_c}$ 使得 $\|\mathbf{W}_c^*\|_F \leq b_{W_c}$ ;  $\mathbf{x}_c = [\mathbf{x}_1^T(t) \ z_1^T(t) \ \mathbf{x}_d^T(t)]^T$ 为评判神经网络的输入量. 基于高斯基函数, 径向基神经网络的 $\phi_c(\mathbf{x}_c)$ 的表达式为

$$\phi_c(\mathbf{x}_c) = \exp\left[-\frac{\|\mathbf{x}_c - \mathbf{c}_r\|}{b_r^2}\right], \quad (23)$$

其中:  $\mathbf{c}_r$ 为径向基神经网络的中心值, 同时,  $b_r$ 为径向基神经网络的宽度. 进一步, 采用时间差分算法, 可得

$$\mathbf{e}_c = \hat{J}(t) - \gamma \hat{J}(t-T) + \mathbf{p}_c = \hat{\mathbf{W}}_c^T \Delta \phi_c(t) + \mathbf{p}_c = \hat{\mathbf{W}}_c^T \Delta \phi_c(t) + \hat{\mathbf{p}}_c + \mathbf{W}_c^{*T} \Delta \phi_c(t), \quad (24)$$

其中 $\Delta \phi_c(t) = [\phi_c(\mathbf{x}_c(t)) - \gamma \phi_c(\mathbf{x}_c(t-T))]$ . 进一步, 可得 $\|\Delta \phi_c(t)\| \leq (1 + \gamma)b_{\phi_c}$ . 最终, 评判网络的权重 $\hat{\mathbf{W}}_c$ 更新律, 设计为

$$\dot{\hat{\mathbf{W}}}_c = -\Gamma_c \Delta \phi_c(t) [\hat{\mathbf{W}}_c^T \Delta \phi_c(t) + \mathbf{p}_c]^T - \sigma_c \Gamma_c \hat{\mathbf{W}}_c, \quad (25)$$

其中:  $\Gamma_c = \text{diag}\{\Gamma_{c1}, \Gamma_{c2}, \Gamma_{c3}\} > 0$ 为神经网络的学习率;  $\sigma_c \Gamma_c \hat{\mathbf{W}}_c$ 为神经网络权重更新率的修正项, 用于放宽神经网络的权重参数学习的持续激励约束. 在式(25)中, 第1项用来减小估计误差 $\mathbf{e}_c = \frac{1}{2} \mathbf{e}_c^T \mathbf{e}_c$ , 同时, 第2项用来增强神经网络权重参数学习过程的鲁棒性, 并解决持续激励缺失引起的评价网络权重更新梯度消失问题. 至此, 强化学习的评判网络设计完成.

### 3.3 融合强化学习执行网络的固定时间控制算法设计

将跟踪误差 $\mathbf{z}_2 = [z_{21} \ z_{22} \ z_{23}]^T$ 关于时间求导数, 可得

$$\dot{\mathbf{z}}_2 = \mathbf{f}_2 + \mathbf{g}_2 \mathbf{N}(\mathbf{v}(t)) + \mathbf{D} - \dot{\boldsymbol{\alpha}}_{2f}, \quad (26)$$

进一步, 可得

$$\mathbf{z}_2^T \dot{\mathbf{z}}_2 = \mathbf{z}_2^T \mathbf{g}_2 \mathbf{N}(\mathbf{v}(t)) + \mathcal{P}_2, \quad (27)$$

其中  $\mathcal{P}_2 = \mathbf{z}_2^T \mathbf{f}_2 - \mathbf{z}_2^T \dot{\boldsymbol{\alpha}}_{2f} + \mathbf{z}_2^T \mathbf{D}$ .

设计如下的Lyapunov函数:

$$\begin{cases} V_2 = V_{21} + V_{22}, V_{22} = \frac{g_2 l_v}{2} \text{tr}(\tilde{\mathbf{W}}_c^T \boldsymbol{\Gamma}_c^{-1} \tilde{\mathbf{W}}_c), \\ V_{21} = V_1 + \frac{1}{2} \mathbf{z}_2^T \mathbf{z}_2 + \frac{g_2}{2 \ell_2} \tilde{b}_2^T \tilde{b}_2 + \\ \frac{g_2}{2} \text{tr}(\tilde{\mathbf{W}}_a^T \boldsymbol{\Gamma}_a^{-1} \tilde{\mathbf{W}}_a), \end{cases} \quad (28)$$

其中:  $\tilde{b}_2 = b_2 - \hat{b}_2$  为参数估计误差,  $\hat{b}_2$  为  $b_2$  的估计值;  $\ell_2 > 0$  为设计参数. 对  $V_{21}$  求关于时间的导数, 可得

$$\begin{aligned} \dot{V}_{21} \leq & \dot{V}_1 + \mathbf{z}_2^T g_2 \mathbf{N}(\mathbf{v}(t)) + \mathcal{P}_2 - \frac{g_2}{\ell_2} \tilde{b}_2^T \dot{\hat{b}}_2 + \\ & g_2 \text{tr}(\tilde{\mathbf{W}}_a^T \boldsymbol{\Gamma}_a^{-1} \dot{\tilde{\mathbf{W}}}_a), \end{aligned} \quad (29)$$

进一步, 再次结合杨氏不等式, 可得

$$\begin{cases} \mathbf{z}_2^T \mathbf{f}_2 \leq g_2 \|\mathbf{z}_2\|^2 \|\mathbf{f}_2\|^2 + \frac{1}{4g_2}, \\ -\mathbf{z}_2^T \dot{\boldsymbol{\alpha}}_{2f} \leq g_2 \|\mathbf{z}_2\|^2 \|\dot{\boldsymbol{\alpha}}_{2f}\|^2 + \frac{1}{4g_2}, \end{cases} \quad (30)$$

进一步, 可得

$$\mathcal{P}_2 \leq g_2 b_2 \|\mathbf{z}_2\|^2 \Pi_2 + \Delta_2 + \mathbf{z}_2^T \mathbf{D}, \quad (31)$$

其中:  $b_2 = \max\{1, \|\mathbf{f}_2\|^2\}$ ,  $\Pi_2 = 1 + \|\boldsymbol{\eta}_1\|^2 \|\mathbf{z}_1\|^2 + \|\dot{\boldsymbol{\alpha}}_{2f}\|^2$ .  $\Delta_2 = \frac{1}{2g_2} > 0$ . 同时, 考虑到电推进器的死

区特性,  $\Delta \mathbf{v}(t)$  为未知量, 以及系统受到的外部有界干扰  $\mathbf{D}$ , 本文引入如下的复合未知项, 具体为<sup>[12, 14]</sup>

$$\boldsymbol{\Xi} = \mathbf{D} + g_2 \Delta \mathbf{v}(t). \quad (32)$$

本文借助于强化学习的Actor网络对上述未知量进行估计, 在消除未知量影响的同时, 优化系统的控制效果, 具体为

$$\boldsymbol{\Xi} = \mathbf{W}_a^* \boldsymbol{\phi}_a(\mathbf{x}_a) + \boldsymbol{\varepsilon}_a(\mathbf{x}_a), \quad (33)$$

其中:  $\mathbf{W}_a^* \in \mathbb{R}^{N_a \times n}$  为理想权重, 满足  $\|\mathbf{W}_a^*\|_F \leq b_{W_a}$ ,  $b_{W_a}$  为大于零的正常数;  $N_a$  为动作网络的隐层;  $\mathbf{x}_a$  为强化学习动作网络的输入.  $\boldsymbol{\phi}_a(\mathbf{x}_a) \in \mathbb{R}^{N_a}$  为径向基神经网络的激活函数, 类似于  $\boldsymbol{\phi}_c$ , 见式(23). 同理, 其满足  $\|\boldsymbol{\phi}_a(\mathbf{x}_a)\| \leq b_{\phi_a}$ , 并且  $b_{\phi_a}$  为大于零的正整数. 由于  $\mathbf{W}_a^*$  为理想权重, 无法直接获得, 在本文中, 采用其估计, 具体如下:

$$\hat{\boldsymbol{\Xi}} = \hat{\mathbf{W}}_a^T \boldsymbol{\phi}_a(\mathbf{x}_a). \quad (34)$$

强化学习动作网络的目标有两个:

1) 使得越界的系统状态能够重新恢复到安全约束界限内, 起到系统状态软约束的目的;

2) 优化系统的控制输出, 使其输出能量越小越好,

即使得评判网络的估计输出  $\hat{\mathbf{J}}(t)$  最小化来逼近去最优值  $\mathbf{J}_d = \mathbf{0}$ .

从而定义强化学习动作网络的估计误差为

$$\mathbf{e}_a = \mathbf{z}_2 + \hat{\mathbf{J}}(t) - \mathbf{J}_d = \mathbf{z}_2 + \hat{\mathbf{W}}_c^T \boldsymbol{\phi}_c(\mathbf{x}_c), \quad (35)$$

由此, 设计  $\hat{\mathbf{W}}_a$  的更新律如下:

$$\dot{\hat{\mathbf{W}}}_a = \boldsymbol{\Gamma}_a \boldsymbol{\phi}_a(t) [\mathbf{z}_2 + \hat{\mathbf{W}}_c^T \Delta \boldsymbol{\phi}_c(\mathbf{x}_c)]^T - \ell_a \boldsymbol{\Gamma}_a \hat{\mathbf{W}}_a, \quad (36)$$

其中:  $\boldsymbol{\Gamma}_a = \text{diag}\{\Gamma_{a1}, \Gamma_{a2}, \Gamma_{a3}\} > 0$  为神经网络的学习率;  $\ell_a \boldsymbol{\Gamma}_a \hat{\mathbf{W}}_a$  为神经网络权重更新率的修正项, 用于放宽神经网络的权重参数学习的持续激励约束. 在式(36)中, 第1项用于第2项的功能与评判网络权重  $\hat{\mathbf{W}}_c$  更新律相同.

进一步, 设计如下的系统的控制器:

$$\begin{cases} \mathbf{v} = -(\mathbf{k}_2 \mathbf{z}_2 + \hat{\mathbf{b}}_2 \mathbf{z}_2 \Pi_2 + l_3 \text{sig}^{\frac{m_2}{n_2}}(\mathbf{z}_2) + \\ \hat{\mathbf{W}}_a^T \boldsymbol{\phi}_a(\mathbf{x}_a) + l_4 \text{sig}^{\frac{p_2}{r_2}}(\mathbf{z}_2)), \\ \dot{\hat{b}}_2 = \ell_2 (\|\mathbf{z}_2\|^2 \Pi_2 - \ell_{\hat{b}_2} \hat{b}_2), \hat{b}_2(0) \geq 0, \end{cases} \quad (37)$$

其中:

$$\begin{aligned} \text{sig}^{\frac{m_2}{n_2}}(\mathbf{z}_2) &= \\ & [\text{sig}^{\frac{m_2}{n_2}}(z_{2,1}) \text{sig}^{\frac{m_2}{n_2}}(z_{2,2}) \text{sig}^{\frac{m_2}{n_2}}(z_{2,3})]^T, \\ \text{sig}^{\frac{p_2}{r_2}}(\mathbf{z}_2) &= \\ & [\text{sig}^{\frac{p_2}{r_2}}(z_{2,1}) \text{sig}^{\frac{p_2}{r_2}}(z_{2,2}) \text{sig}^{\frac{p_2}{r_2}}(z_{2,3})]^T. \end{aligned}$$

## 4 稳定性证明

**定理1** 考虑带有死区特性电推进器的航天器系统(4)–(6), 针对其状态安全约束控制(7), 采用固定时间控制器及其自适应律(16)(37), 结合强化学习评判网络(18)及其更新律(25), 以及用于克服死区约束的控制输出与越界状态的强化学习动作网络(33)及其更新律(36), 可得以下结论:

1) 所有闭环系统的信号都是一致最终有界的. 其收敛域为

$$\begin{aligned} \|\tilde{\mathbf{W}}_a\|_F &\leq \sqrt{\frac{2(\mathcal{Y} + V_2(0))}{g_2 \lambda_{\min}(\boldsymbol{\Gamma}_a^{-1})}}, \\ \|\tilde{\mathbf{W}}_c\|_F &\leq \sqrt{\frac{2(\mathcal{Y} + V_2(0))}{l_v \lambda_{\min}(\boldsymbol{\Gamma}_c^{-1})}}, \\ \|\tilde{b}_i\|_F &\leq \sqrt{\frac{2\gamma_i(\mathcal{Y} + V_2(0))}{g_i}}, \quad i = 1, 2, \\ \|\mathbf{y}_2\|_F &\leq \sqrt{2(\mathcal{Y} + V_2(0))}. \end{aligned}$$

2) 系统状态可在固定时间收敛至以下约束区域:

$$\begin{aligned} \|z_{1,j}\| &\leq \min\left\{\left(\frac{M}{(1-\theta)a}\right)^{\frac{1}{n+m}}, \left(\frac{M}{(1-\theta)b}\right)^{\frac{1}{r+p}}\right\}, \\ t_{s1} &\leq \frac{1}{a} \frac{2n}{m-n} + \frac{1}{b} \frac{2r}{r-p}, \end{aligned}$$

可通过适当参数调节来调节系统状态的收敛时间及收敛域。

3) 系统的状态将会保持在式(7)定义的安全约束范围内运行, 一旦受扰, 偏离安全约束范围, 强化学习惩罚机制(18)将使其重新恢复到约束域中。

证 结合设计的 Lyapunov 函数, 并代入控制律(37)及其自适应律, 可得

$$\begin{aligned} \dot{V}_{21} \leq & -\bar{\varepsilon}_1 \mathbf{y}_1^T \mathbf{y}_1 - \sum_{i=1}^2 g_i \lambda_{\min}(\mathbf{k}_i) \|\mathbf{z}_i\|^2 - \\ & \sum_{i=1}^2 \frac{g_i \ell_{\tilde{b}_i}}{2} \|\tilde{b}_i\|^2 + \sum_{i=1}^2 \Gamma_i + \tilde{\mathbf{h}}_1^T \tilde{\mathbf{h}}_1 - \\ & g_2 \mathbf{z}_2^T \tilde{\mathbf{W}}_a^T \phi_a(\mathbf{x}_a) + \mathbf{z}_2^T g_2 \varepsilon_a(\mathbf{x}_a) - \\ & l_a g_2 \text{tr}(\tilde{\mathbf{W}}_a^T \hat{\mathbf{W}}_a) + g_2 \text{tr}(\tilde{\mathbf{W}}_a^T \Delta \phi_a(t) \mathbf{z}_2^T) + \\ & g_2 \text{tr}(\tilde{\mathbf{W}}_a^T \Delta \phi_a(t) [\hat{\mathbf{W}}_c^T \Delta \phi_c(\mathbf{x}_c)]^T), \end{aligned} \quad (38)$$

其中  $\Gamma_2 = \Delta_2 + \frac{\ell_{\tilde{b}_2}}{2} \|\mathbf{b}_2\|^2$ . 进一步, 结合不等式

$$\begin{aligned} l_a g_2 \text{tr}(\tilde{\mathbf{W}}_a^T \hat{\mathbf{W}}_a) & \leq -\frac{l_a g_2}{2} \text{tr}(\tilde{\mathbf{W}}_a^T \tilde{\mathbf{W}}_a) + \frac{l_a g_2 b_{W_a}^2}{2}, \\ -\mathbf{z}_2^T g_2 k_{21} \mathbf{z}_2 + \mathbf{z}_2 g_2 \varepsilon_a(\mathbf{x}_a) & \leq \\ -g_2 \lambda_{\min}(\mathbf{k}_{21}) \|\mathbf{z}_2\|^2 + g_2 b_{\varepsilon_a} \|\mathbf{z}_2\| & \leq \frac{g_2 b_{\varepsilon_a}^2}{4 \lambda_{\min}(\mathbf{k}_{21})}, \\ -l_a \|\tilde{\mathbf{W}}_a\|_F^2 + b_{\phi_a} b_{\phi_c} b_{W_c} \|\tilde{\mathbf{W}}_a\|_F & \leq \frac{b_{\phi_a}^2 b_{\phi_c}^2 b_{W_c}^2}{4 l_a}, \\ b_{\phi_a} b_{\phi_c} \|\tilde{\mathbf{W}}_a\|_F \|\tilde{\mathbf{W}}_c\|_F & \leq \\ \frac{b_{\phi_a} b_{\phi_c}}{2} \|\tilde{\mathbf{W}}_a\|_F^2 + \frac{b_{\phi_a} b_{\phi_c}}{2} \|\tilde{\mathbf{W}}_c\|_F^2 & \|\hat{\mathbf{W}}_c\|_F \leq \\ \|\tilde{\mathbf{W}}_c\|_F + b_{W_c}, \end{aligned} \quad (39)$$

进一步, 可得

$$\begin{aligned} \dot{V}_{21} \leq & -g_1 \lambda_{\min}(\mathbf{k}_1) \|\mathbf{z}_1\|^2 - g_2 \lambda_{\min}(\mathbf{k}_{22}) \|\mathbf{z}_2\|^2 - \\ & g_2 \left( \frac{l_a}{2} - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_a \right) \text{tr}(\tilde{\mathbf{W}}_a^T \tilde{\mathbf{W}}_a) - \bar{\varepsilon}_1 \mathbf{y}_1^T \mathbf{y}_1 + \\ & g_2 \frac{b_{\phi_a} b_{\phi_c}}{2} \text{tr}(\tilde{\mathbf{W}}_c^T \tilde{\mathbf{W}}_c) - \sum_{i=1}^2 \frac{g_i \ell_{\tilde{b}_i}}{2} \|\tilde{b}_i\|^2 + b_{V_1}, \end{aligned} \quad (40)$$

其中  $b_{V_1} = \frac{l_a g_2 b_{W_a}^2}{2} + \frac{g_2 b_{\phi_a}^2 b_{\phi_c}^2 b_{W_c}^2}{4 l_a} + \sum_{k=1}^2 \Gamma_k + \tilde{\mathbf{h}}^T \tilde{\mathbf{h}}$ .

进一步, 将  $V_{22}$  关于时间求导数, 可得

$$\dot{V}_{22} \leq -g_2 l_v \sigma_c \text{tr}(\tilde{\mathbf{W}}_c^T \tilde{\mathbf{W}}_c) + b_{V_2} (\|\tilde{\mathbf{W}}_c\|_F), \quad (41)$$

其中  $b_{V_2} = g_2(1 + \gamma) l_v b_{\phi_c} b_{p_c} + g_2(1 + \gamma)^2 l_v b_{\phi_c}^2 b_{W_c}$ .

进一步, 可得

$$\begin{aligned} \dot{V}_2 = \dot{V}_{21} + \dot{V}_{22} \leq & -g_1 \lambda_{\min}(\mathbf{k}_1) \|\mathbf{z}_1\|^2 - g_2 \lambda_{\min}(\mathbf{k}_{22}) \|\mathbf{z}_2\|^2 - \\ & g_2 \left( \frac{l_a}{2} - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_a \right) \text{tr}(\tilde{\mathbf{W}}_a^T \tilde{\mathbf{W}}_a) - \bar{\varepsilon}_1 \mathbf{y}_1^T \mathbf{y}_1 + b_{V_1} - \end{aligned}$$

$$\begin{aligned} & g_2 \left( l_v \sigma_c - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_c \right) \text{tr}(\tilde{\mathbf{W}}_c^T \tilde{\mathbf{W}}_c) + b_{V_2} \|\tilde{\mathbf{W}}_c\|_F - \\ & l_c \text{tr}(\tilde{\mathbf{W}}_c^T \tilde{\mathbf{W}}_c) - \sum_{i=1}^2 \frac{g_i \ell_{\tilde{b}_i}}{2} \|\tilde{b}_i\|^2, \end{aligned} \quad (42)$$

结合不等式  $-l_c \|\tilde{\mathbf{W}}_c\|_F^2 + b_{V_2} \|\tilde{\mathbf{W}}_c\|_F \leq \frac{b_{V_2}^2}{2 l_c}$ , 进一步, 可得

$$\dot{V}_2 \leq -\Re V_2 + \Theta, \quad (43)$$

其中:

$$\begin{aligned} \Re = \min \{ & g_1 \lambda_{\min}(\mathbf{k}_1), g_2 \lambda_{\min}(\mathbf{k}_{22}), \\ & \frac{\sigma_a - b_{\phi_a} b_{\phi_c} - 2 l_a}{\lambda_{\min}(\mathbf{\Gamma}_a^{-1})}, \ell_{\tilde{b}_1}, \ell_{\tilde{b}_2}, \bar{\varepsilon}_2, \\ & \frac{2 l_v \sigma_c - b_{\phi_a} b_{\phi_c} - 2 l_c}{l_v \lambda_{\min}(\mathbf{\Gamma}_c^{-1})} \}, \\ \Theta = & b_{V_1} + \frac{b_{V_2}^2}{2 l_c}, \end{aligned}$$

同时对不等式(43)两侧乘以  $e^{\Re t}$ , 并定义  $\Upsilon = \Theta/\Re$ , 在区间  $[0, t]$  上积分式(43)可得

$$V_2(t) \leq \Upsilon + V_2(0) e^{-\Re t}, \quad (44)$$

其中  $V_2(0) = \frac{1}{2} \sum_{i=1}^2 [\mathbf{z}_i^T(0) \mathbf{z}_i(0) + \frac{g_i}{2 \gamma_i} \tilde{b}_i^2(0)] + \frac{1}{2} \times \mathbf{y}_2^T(0) \mathbf{y}_2(0) + \frac{g_2}{2} \text{tr}(\tilde{\mathbf{W}}_a^T(0) \mathbf{\Gamma}_a^{-1} \tilde{\mathbf{W}}_a(0)) + \frac{g_2 l_v}{2} \times \text{tr}(\tilde{\mathbf{W}}_c^T(0) \mathbf{\Gamma}_c^{-1} \tilde{\mathbf{W}}_c(0))$ . 此时, 参数  $l_a, l_v \sigma_c$  的选取满足  $\frac{l_a}{2} - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_a > 0, l_v \sigma_c - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_c > 0$ , 可得  $\mathbf{z}_1, \mathbf{z}_2, \tilde{b}_1, \tilde{b}_2, \mathbf{y}_1, \tilde{\mathbf{W}}_a, \tilde{\mathbf{W}}_c$  都是有界的, 完成定理1中第1部分的证明。

此时, 注意到在紧集  $\mathcal{D}_d = \{\mathbf{y}_d \in \mathbb{R}^3 : -\mathbf{A}_i \mathbf{0} \leq \mathbf{y}_{id} \leq \mathbf{A}_i \mathbf{0}, i = 1, 2, 3\}$  中,  $\alpha_0 \in L_\infty$ , 结合  $\mathbf{z}_1 = \zeta_1 - \alpha_0$ , 可知  $\zeta_1$  为有界的. 由于有界的

$$\frac{(F_{i1} F_{i2} + x_{1,i}^2)}{[(F_{i1} + x_{1,i})(F_{i2} - x_{1,i})]^2},$$

可得  $\Pi_1 \in L_\infty$ , 进一步结合式(16)可得  $\alpha_1 \in L_\infty, \dot{b}_1 \in L_\infty$ . 随后, 结合式(9), 可得  $\alpha_{2f}$  为有界量. 同理, 可分析  $\zeta_2, \alpha_2, \dot{b}_2$  为有界量, 以及系统的控制输入  $\mathbf{v}$  也是有界量. 同时, 根据  $b_1, b_2$  的有界定义, 可得  $\hat{b}_1, \hat{b}_2$  也是有界量. 再次基于  $\mathbf{z}_1 = \zeta_1 - \alpha(0)$ , 结合式(7), 可得

$$\delta_1 \mathbf{z}_1 = \delta_2 \mathbf{e}, \quad (45)$$

其中: 当  $-F_{i1} < x_{1i} < F_{i2}$  时, 可得  $\delta_1 = [\delta_{1,1} \ \delta_{1,2} \ \delta_{1,3}]^T$ , 具体为  $\delta_{1i} = (F_{i1} - x_{1i})(F_{i2} + x_{1i})(F_{i1} - y_{di})(F_{i2} + y_{di}) > 0$ , 同时,  $\delta_2 = [\delta_{2,1} \ \delta_{2,2} \ \delta_{2,3}]^T$ , 具体是  $\delta_{2i} = F_{i1} F_{i2} + x_{1i} y_{di} > 0$ , 此时, 跟踪误差满足  $\mathbf{e} = \mathbf{x}_1 - \mathbf{y}_d$ . 进一步

$$\mathbf{e} = \delta \mathbf{z}_1, \quad (46)$$

其中  $\delta = \text{diag}\{\frac{\delta_{11}}{\delta_{21}}, \frac{\delta_{12}}{\delta_{22}}, \frac{\delta_{13}}{\delta_{23}}\}$ , 由于  $z_1$  是有界的, 则  $e \in L_\infty$ . 至此, 所有的闭环系统的信号都是有界的. 进一步, 基于上述有界性的分析, 可知, 存在正常数  $M$  使得  $-\sum_{i=1}^2 \frac{g_i \ell_{b2}}{2} \|\tilde{b}_i\|^2 - g_2 (\frac{\ell_a}{2} - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_a) \text{tr}(\tilde{W}_a^T \times \tilde{W}_a) - g_2 (l_v \sigma_c - \frac{b_{\phi_a} b_{\phi_c}}{2} - l_c) \text{tr}(\tilde{W}_c^T \tilde{W}_c) - \bar{\varepsilon}_1 \mathbf{y}_1^T \mathbf{y}_1 + \Theta \leq M$ . 在此基础上, 考虑Lyapunov函数  $V_3 = \frac{z_1^T z_1}{2} + \frac{z_2^T z_2}{2}$ , 并对其求关于时间的导数, 可得

$$\dot{V}_3 \leq -aV_3^{\frac{n+m}{2n}} - bV_3^{\frac{r+p}{2r}} + M, \quad (47)$$

其中:  $a = 2^{1-\frac{n+m}{2n}} 2^{\frac{n+m}{2n}} a_1$ ,  $a_1 = \min\{g_1 l_1, g_2 l_3\}$ ,  $\frac{m}{n} = \min\{\frac{m_1}{n_1}, \frac{m_2}{n_2}\}$ ,  $b = 2^{\frac{r+p}{2r}} b_1$ ,  $b_1 = \min\{g_1 l_2, g_2 l_5\}$ ,  $\frac{p}{r} = \max\{\frac{p_1}{r_1}, \frac{p_2}{r_2}\}$ . 同时,  $\frac{m+n}{2n} > 1$ ,  $0 < \frac{p+r}{2r} < 1$ .

结合文献[15]中的引理2, 系统固定时间的收敛域  $\Sigma$  可表示为  $\Sigma = \{\lim_{t \rightarrow t_{s1}} z_1 | V_3 \leq \min\{(\frac{M}{(1-\theta)a})^{\frac{2n}{n+m}}, (\frac{M}{(1-\theta)b})^{\frac{2r}{r+p}}\}\}$ . 系统收敛的固定时间上界为  $t_{s1} \leq \frac{1}{a} \frac{1}{\frac{m+n}{2n} - 1} + \frac{1}{b} \frac{1}{1 - \frac{p+r}{2r}} \leq \frac{1}{a} \frac{2n}{m-n} + \frac{1}{b} \frac{2r}{r-p}$ . 此时, 一旦系统的状态  $z_1$  收敛至  $\Sigma$  内, 在固定时间上界  $t_{s1}$ , 系统状态的收敛域可为

$$|z_{1,j}| \leq \min\{(\frac{M}{(1-\theta)a})^{\frac{n}{n+m}}, (\frac{M}{(1-\theta)b})^{\frac{r}{r+p}}\},$$

至此, 系统的稳定性证明结束. 证毕.

## 5 结果与讨论

### 5.1 仿真初始条件及控制器参数

本节在 MATLAB/Simulink 环境下进行数值仿真验证, 控制目标是考虑外部干扰和执行机构死区特性的条件下, 使航天器编队在 100 s 内达到期望构型, 并满足 0.01 m 的相对位置保持精度和防碰撞要求(两颗星相对距离不小于 1.5 m), 在系统的初始参数如表 1 所示. 本文的仿真环境是 Intel(R) Core(TM) i7-8550U CPU @ 1.80 GHz, 8 G 内存.

表 1 系统仿真参数

Table 1 System parameters for simulation

序号	系统参数	值
1	$m_1$	59 kg
2	$m_f$	59 kg
3	$R_1$	7578.17 km
4	$\omega_0$	$9.57 \times 10^{-4}$ rad/s
5	$R_0(0)$	$[6 \ 7 \ 5]^T$ m
6	$v_0(0)$	$[0 \ 0 \ 0]^T$ m/s

在数学仿真中, 假设两颗星的期望相对位置为  $[3 \ 0 \ 0]^T$  m, 期望相对速度为  $[0 \ 0 \ 0]^T$  (m/s). 在实际应用中, 本文考虑超紧密航天器在近地轨道附近受到的气动阻力、太阳光压等干扰力幅值不超过  $10^{-3}$  N, 故而, 本文假设系统干扰力上限为  $1.5 \times 10^{-3}$  N. 电推进器死区特性的参数为  $h_r = 1, h_l = 1, o_r = 0.3, o_l = -0.3$ .

考虑到系统的状态约束问题, 本文设计的状态约束函数的参数为  $F_{11} = 2, F_{21} = F_{31} = -4, F_{12} = F_{22} = F_{32} = 15$ . 另外, 采用的电推进器死区特性的参数为  $h_r = 1, h_l = 1, o_r = 0.3, o_l = -0.3$ .

根据系统的控制需求, 选取位置控制回路的参数为  $\mathbf{k}_1 = [0.05 \ 0.05 \ 0.05]^T, l_1 = l_2 = 0.01, m_1 = 9, n_1 = 5, p_1 = 5, r_1 = 7, \ell_{b1} = \ell_{b2} = 1$ . 速度控制回路参数为  $\mathbf{k}_2 = [15 \ 15 \ 15]^T, l_3 = 0.1, l_4 = 0.01, m_2 = 11, n_2 = 9, p_2 = 5, r_2 = 7$ .

强化学习评判网络的初始参数为

$$\Gamma_c = 0.1, \sigma_c = 0.2,$$

$$c_{\text{critic}} = \begin{bmatrix} -1.5 & -1 & -0.5 & 0 & 0.5 & 1 & 1.5 \\ -1.5 & -1 & -0.5 & 0 & 0.5 & 1 & 1.5 \\ -1.5 & -1 & -0.5 & 0 & 0.5 & 1 & 1.5 \end{bmatrix},$$

$$b_{\text{critic}} = 2,$$

此外, 强化学习动作网络的初始参数为

$$\Gamma_a = 0.01, \sigma_a = 0.1, c_{\text{actor}} = c_{\text{critic}}, b_{\text{actor}} = 2.$$

其余参数为  $R_s = I_3, Q_s = I_7, I_n$  为  $n$  维的单位矩阵.  $\varepsilon = 0.1, \varphi = 10$ . 仿真时间为 400 s, 仿真结果如图 3-10 所示.

图 3 与图 4 分别为超紧密编队航天器的三维相对位置曲线和相对位置随时间变化曲线, 图 5 为编队航天器的相对速度随时间变化的曲线. 从图 3-5 可以看出, 两颗星在 90 s 达到期望的相对位置, 满足固定时间的控制要求, 相对位置保持精度达到 0.031 m.

图 6 为强化学习奖励函数变化曲线, 图 7 和图 8 分别为强化学习评判网络中的函数  $\hat{W}_c^T \Delta \phi_c$  和强化学习动作网络的补偿量  $\hat{W}_a^T \Delta \phi_a$  变化曲线. 从图中可以看出在航天器编队相对位置重构的过程中, 在 20 s 左右, 超紧密编队行为出现的状态越界的行为, 此时, 提出的基于强化学习的软约束机制开始工作: 由图 6 可以看出, 此时出现相对位置轴方向的越界行为, 致使轴方向的强化学习奖励函数值急剧增大, 与目标奖励函数越小越好 ( $p_{c1}$  越小越好) 的相背离. 进而随着  $p_{c1}$  的增大, 强化学习评判网络中的  $(\hat{W}_c^T \Delta \phi_c)_1$  (见于图 7) 也迅速增大, 进一步引发图 6 中强化学习动作网络的在  $x$  轴方向的补偿量  $(\hat{W}_a^T \Delta \phi_a)_1$  也迅速增大, 形成负反馈调节机制, 调节超紧密编队航天器  $x$  轴方向的运动位置  $\rho_1$  (见图 4), 重新回归到安全的避碰约束范围 (1.5 m) 外以满足超紧密编队的防碰撞的控制需求.



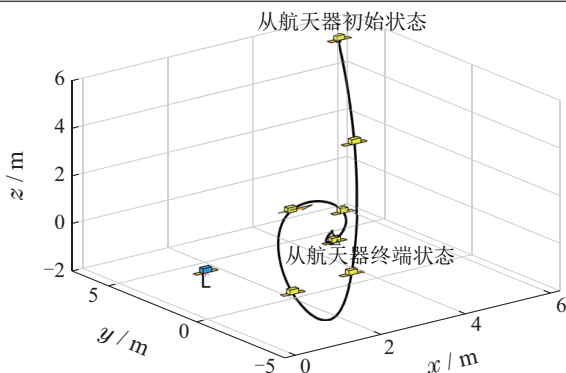


图3 航天器超紧密编队三维运行轨迹

Fig. 3 Moving trajectories of the spacecraft in ultra close formation in view of 3-D

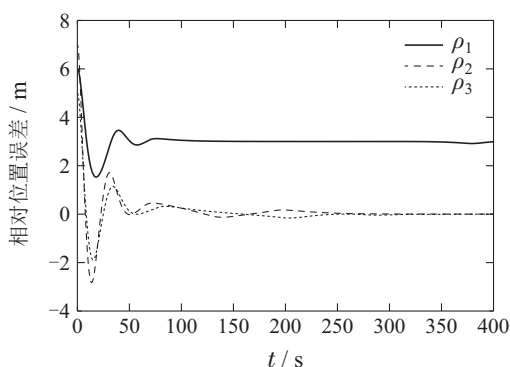


图4 航天器超紧密编队的相对位置误差

Fig. 4 Trajectories of relative position error

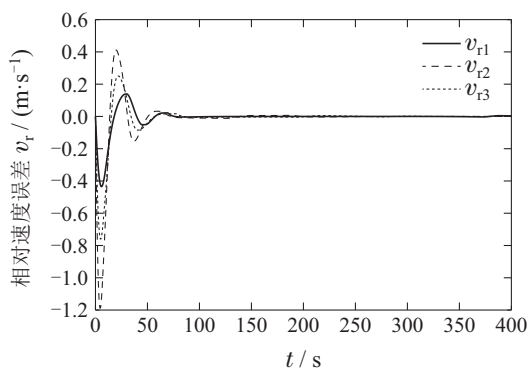


图5 航天器超紧密编队相对速度误差

Fig. 5 Trajectories of relative velocity error

为进一步说明所提出的算法能减小执行机构存在的死区带来的精度影响, 本节在控制器设计中, 不考虑基于强化学习的死区补偿机制, 在同样的仿真参数下, 仿真结果如图9-10所示. 图9与图10分别为超紧密编队航天器的三维相对位置和相对距离随时间变化的曲线. 由于减小外部扰动的影响, 推力器在构型保持阶段一直在工作, 以保持超紧密编队的高精度控制行为. 由于执行机构死区效应的影响, 致使执行机构无法实时响应超紧密编队高精度的控制需求, 如图10所示, 如不进行基于强化学习的实施补偿, 航天器相对位置不能满足保持0.01 m的相对位置精度的控制需

求, 最终只能达到0.3 m左右的位置保持精度, 对比图4, 可以看出, 本文所提出的基于强化学习的死区补偿机制能够处理死区效应带来的控制精度影响, 可以使航天器编队的相对位置保持精度达到0.01 m. 仿真结果表明本文所提出控制方法的有效性.

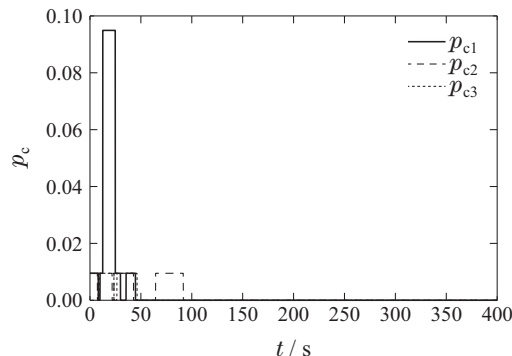


图6 强化学习奖励函数曲线

Fig. 6 Reinforcement learning reward function curves

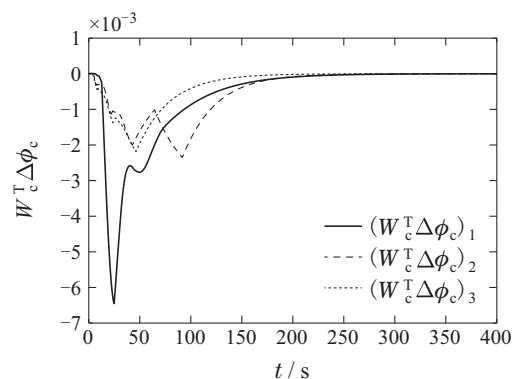


图7 强化学习评判网络输出曲线

Fig. 7 Reinforcement learning critic network function

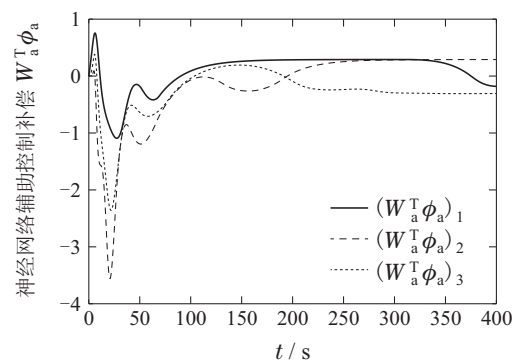


图8 强化学习动作网络对死区效应的补偿量

Fig. 8 Dead-zone compensation output via action network of reinforcement learning

## 6 结论

本文对近地轨道超紧密航天器的编队重构问题进行研究, 考虑到电推进器的死区特性、系统受到的外界环境干扰以及避免碰撞的安全环约束需求, 设计了一种基于强化学习的电推力固定时间超紧密编队控制器设计方法: 1) 将超紧密编队构型保持转化为依赖

于系统状态的状态约束控制,同时,考虑到紧密编队建立构型的时间约束,引入固定时间的控制设计,保障系统在外外部扰动作用下,能够在有限时间内到达满足编队构型约束的状态约束控制目标;2)考虑到电推进器死区效应的影响,本文引入强化学习评判-动作机制,通过最小化评判网络输出的方式,辅助以动作网络的辅助补偿,来解决死区特性带来的影响.与此同时,在该评判-动作控制架构上设计状态越界软约束拉回机制,解决因外部扰动引起的航天器碰撞与紧密编队构型破坏问题,提高航天器紧密编队的鲁棒性;3)在后续的研究中,将在此基础上,研究多航天器超紧密编队的分布式协同优化控制问题,设计满足全局优化的协同编队控制算法.

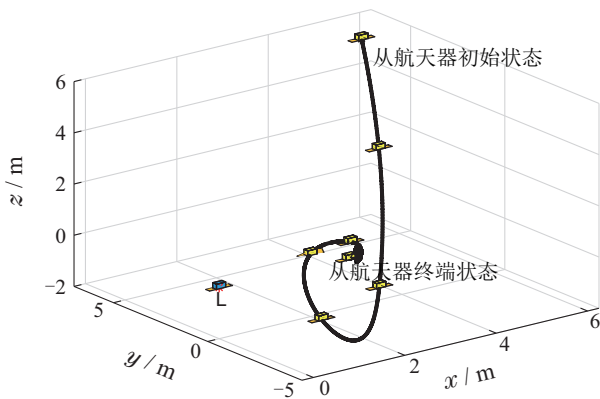


图9 忽略无死区效应的航天器超紧密编队三维运行轨迹  
Fig. 9 Moving trajectories of the spacecraft formation without dead-zone compensation in view of 3-D

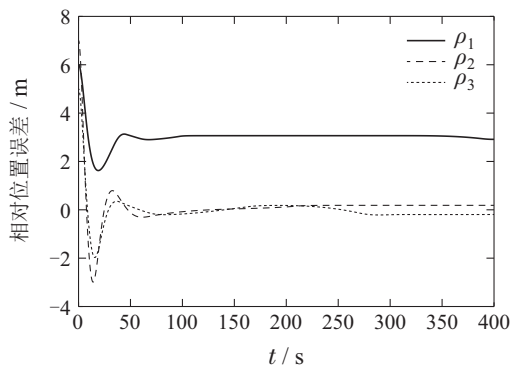


图10 忽略无死区效应的航天器超紧密编队相对位置误差  
Fig. 10 Trajectories of relative position without function of dead-zone compensation

## 参考文献:

- [1] HUANG Jing, SUN Lujun, SUN Jun, et al. Cooperative control for collision avoidance of ultra-close spacecraft formation based on prescribed performance control. *Flight Control and Detection*, 2019, 2(3): 52 – 60.  
(黄静, 孙禄君, 孙俊, 等. 基于预设性能控制的超紧密航天器编队防撞协同控制. 飞控与探测, 2019, 2(3): 52 – 60.)
- [2] MAYNE D Q, RAWLINGS J B, RAO C V, et al. Constrained model predictive control: Stability and optimality. *Automatica*, 2000, 36(6): 789 – 814.
- [3] BURGER M, GUAY M. Robust constraint satisfaction for continuous time nonlinear systems in strict feedback form. *IEEE Transactions on Automatic Control*, 2010, 55(11): 2597 – 2601.
- [4] NGO K B, MAHONY R, JIANG Z P. Integrator backstepping using barrier functions for systems with multiple state constraints. *Proceedings of the 44th IEEE Conference on Decision and Control*. Seville, Spain: IEEE, 2005: 8306 – 8312.
- [5] HU Q, SHAO X, GUO L. Adaptive fault-tolerant attitude tracking control of spacecraft with prescribed performance. *IEEE/ASME Transactions on Mechatronics*, 2018, 23(1): 331 – 341.
- [6] ZHAO K, SONG Y D. Removing the feasibility conditions imposed on tracking control designs for state-constrained strict-feedback systems. *IEEE Transactions on Automatic Control*, 2019, 64(3): 1265 – 1272.
- [7] AN L, YANG G H. Collisions-free distributed optimal coordination for multiple Euler-Lagrangian systems. *IEEE Transactions on Automatic Control*, 2021, 67(1): 460 – 467.
- [8] HONG Huifen. *Research on finite-time consensus of multi-agent systems with limited communication and its related problems*. Nanjing: Southeast University, 2019.  
(洪会粉. 通信受限的多智能体系统有限时间一致性及其相关问题研究. 南京: 东南大学, 2019.)
- [9] HAN Luofeng, ZHU Kangwu, HUANG Wenbin, et al. Development and on orbit verification of M2 microwave ion propulsion system. *Vacuum and Cryogenics*, 2022, 28(1): 98 – 105.  
(韩罗峰, 朱康武, 黄文斌, 等. M2型微波离子推进系统研制及在轨验证. 真空与低温, 2022, 28(1): 98 – 105.)
- [10] LV Y, HU Q, MA G, et al. 6 DOF synchronized control for spacecraft formation flying with input constraint and parameter uncertainties. *ISA Transactions*, 2011, 50(4): 573 – 580.
- [11] ZHANG D, KONG L, ZHANG S, et al. Neural networks-based fixed-time control for a robot with uncertainties and input deadzone. *Neurocomputing*, 2020, 390: 139 – 147.
- [12] HE W, HUANG B, DONG Y T, et al. Adaptive neural network control for robotic manipulators with unknown deadzone. *IEEE Transactions on Cybernetics*, 2018, 48(9): 2670 – 2682.
- [13] GUO X, YAN W, CUI R. Integral reinforcement learning-based adaptive NN control for continuous-time nonlinear MIMO systems with unknown control directions. *IEEE Transactions on Systems, Man, and Cybernetics*, 2020, 50(11): 4068 – 4077.
- [14] LI H, WU Y, CHEN M. Adaptive fault-tolerant tracking control for discrete-time multiagent systems via reinforcement learning algorithm. *IEEE Transactions on Cybernetics*, 2021, 51(3): 1163 – 1174.
- [15] CHEN Q, XIE S, SUN M, et al. Adaptive nonsingular fixed-time attitude stabilization of uncertain spacecraft. *IEEE Transactions on Aerospace and Electronic Systems*, 2018, 54(6): 2937 – 2950.

## 作者简介:

孟亦真 中级工程师, 目前研究方向为航天器集群智能避障与逃逸博弈控制, E-mail: 15151830168@163.com;

黄静 高级工程师, 目前研究方向为航天器振动抑制与容错控制, E-mail: huangjing04415@163.com;

周绍辉 高级工程师, 目前研究方向为航天器电离子推进器控制, E-mail: zsh@sastspace.com;

周彬 中级工程师, 目前研究方向为航天器姿态最优控制, E-mail: 732200801@qq.com;

朱康武 副研究员, 目前研究方向为航天器电离子推进器结构设计, E-mail: zjuzkw@zju.edu.cn.