

控制与决策

Control and Decision

基于增量式Q学习的固定翼无人机跟踪控制性能优化

赵振根, 程磊

引用本文:

赵振根,程磊. 基于增量式Q学习的固定翼无人机跟踪控制性能优化[J]. *控制与决策*, 2024, 39(2): 391–400.

在线阅读 View online: <https://doi.org/10.13195/j.kzyjc.2022.0708>

您可能感兴趣的其他文章

Articles you may be interested in

基于领航-跟随的有人/无人机编队队形保持控制

Formation keeping control for manned/unmanned aerial vehicle formation based on leader-follower strategy

控制与决策. 2021, 36(10): 2435–2441 <https://doi.org/10.13195/j.kzyjc.2020.0453>

输出误差约束下四旋翼无人机预定性能反步控制

Prescribed performance backstepping control for quadrotor UAV with output error constraint

控制与决策. 2021, 36(5): 1059–1068 <https://doi.org/10.13195/j.kzyjc.2019.1249>

分布式无人机的时变编队非线性控制设计

Time-varying formation nonlinear control of distributed multiple UAVs

控制与决策. 2021, 36(10): 2490–2496 <https://doi.org/10.13195/j.kzyjc.2020.0136>

基于数据驱动的非线性网络系统自适应迭代学习控制

Data driven adaptive learning control of nonlinear network system

控制与决策. 2021, 36(6): 1523–1528 <https://doi.org/10.13195/j.kzyjc.2019.1182>

参数未知的离散系统Q-学习优化状态估计与控制

Q-learning optimal state estimation and control for discrete systems with unknown parameters

控制与决策. 2020, 35(12): 2889–2897 <https://doi.org/10.13195/j.kzyjc.2019.0180>

基于增量式 Q 学习的固定翼无人机跟踪控制性能优化

赵振根[†], 程磊

(南京航空航天大学 自动化学院, 南京 211106)

摘要: 针对固定翼无人机纵向控制的高性能需求, 提出一种控制系统性能优化结构. 该结构包括一个使系统稳定的标称控制器和一个参与性能优化的增量式控制器. 控制系统增量式的实现不会改变原有的控制系统, 而是仅对标称控制系统做控制输入的补偿与控制性能的优化. 基于 Q 学习理论进行增量式控制器设计, 针对状态信息完全可获得的系统, 设计一种基于状态反馈的增量式 Q 学习算法. 当状态信息不能完全获得时, 利用系统输入、输出和参考信号数据, 设计一种基于输出反馈的增量式 Q 学习算法. 两种增量式控制器均是在数据驱动环境下自适应学习增量式控制律, 无需提前知道系统动力学模型以及标称控制器的控制增益. 此外, 证明了增量式 Q 学习方法在满足持续激励条件的激励噪声下, 对 Q 函数贝尔曼方程的求解没有偏差. 最后, 通过对 F-16 飞行器纵向模型实例的仿真验证该方法的有效性.

关键词: 强化学习; Q 学习; 增量式控制; 性能优化; 跟踪控制; 无人机

中图分类号: TP273

文献标志码: A

DOI: 10.13195/j.kzyjc.2022.0708

引用格式: 赵振根, 程磊. 基于增量式 Q 学习的固定翼无人机跟踪控制性能优化[J]. 控制与决策, 2024, 39(2): 391-400.

Performance optimization for tracking control of fixed-wing UAV with incremental Q -learning

ZHAO Zhen-gen[†], CHENG Lei

(College of Automation Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China)

Abstract: Aiming at the high performance requirements of longitudinal control of a fixed-wing unmanned aerial vehicle (UAV), a performance optimization structure of the control system is proposed. This structure includes a nominal controller that stabilizes the system and an incremental controller that participates in performance optimization. The incremental implementation of the control system does not change the original control system, but compensates the control input and optimizes the control performance for the nominal control system exclusively. Based on the Q -learning theory, the incremental controller is designed. For the system with completely available state information, an incremental Q -learning algorithm based on state feedback is developed. When the state information cannot be obtained completely, an incremental Q -learning algorithm based on output feedback is designed by using the system input, output and reference trajectory data. Both incremental controllers learn incremental control laws adaptively in the data-driven environment without the need for system dynamics model and the control gain of the nominal controller. In addition, it is proved that the incremental Q -learning method has no bias in solving the Q -function Bellman equation under the excitation noise. Finally, the effectiveness of the method is verified by the simulation of an example of the longitudinal model of the F-16 aircraft.

Keywords: reinforcement learning; Q -learning; incremental control; performance optimization; tracking control; UAV

0 引言

近些年来, 无人机逐步成为当今航空航天领域的研究热点. 无人机的应用领域广泛、型号多样, 在人类生活

生产中都得到了极大的应用与普及, 例如航拍摄影、无人投递、侦察探测等. 但随着无人机的任务多样化、结构复杂化, 无人机也面临越来越多的挑战. 飞行控

收稿日期: 2022-04-27; 录用日期: 2022-10-10.

基金项目: 国家自然科学基金项目(62003161); 江苏省自然科学基金项目(BK20190399); 中国博士后科学基金项目(2021M701701).

责任编辑: 侯忠生.

[†]通讯作者. E-mail: zhaozhengen@nuaa.edu.cn.

*本文附带电子附录文件, 可登录本刊官网该文“资源附件”区自行下载阅览.

制系统作为无人机的“大脑”,不仅需要提高对环境的适应性,还要能够处理自身复杂的动力学模型和自身模型的不确定性,从而使无人机表现出一定的自主性与智能性.目前,飞行控制系统已经得到了广泛的研究,并取得了许多成果.

传统控制方法发展已相对成熟,例如PID控制、切换控制、动态逆控制、鲁棒控制和自适应控制等^[1].文献[2]针对无人机飞行过程中出现的参数摄动和外部扰动的影响,提出在经典PID控制的基础上加入 H_∞ 鲁棒控制的混合控制方式.文献[3]将非线性动态逆控制与反步法相结合,对高超音速飞机纵向运动模型设计飞行控制系统.该系统以非线性动态逆控制作为控制内环,以反步法作为控制外环从而保证系统的全局稳定以及抑制不确定参数的扰动.

传统无人机控制设计方法大多依赖于系统动力学模型,无模型的PID控制的参数整定一般需要用试凑法反复调试才能实现,存在调参繁琐等问题.强化学习对无模型控制问题提供了一种解决思路.文献[4]提出了一种利用强化学习技术训练的神经网络来控制四旋翼的方法,所使用策略网络对阶跃响应能够相对准确地做出反应.文献[5]将基于表格的Q学习应用于实际的四旋翼无人机以实现悬停控制.文献[6]提出了一种基于深度强化学习算法的鲁棒控制方案,该方案利用神经网络实现了状态到控制指令的端对端设计.对于无人机非线性系统的处理,代尔夫特理工大学研究团队提供了一些思路,文献[7-8]利用增量式近似动态规划处理非线性系统,它结合了线性近似动态规划方法和增量控制技术,通过系统辨识得到系统的增量模型,并在线设计最优控制器,最后对导弹俯仰平面,某固定翼飞机的高度控制以及航天器姿态控制的仿真验证了该方法的有效性.

上述关于无人机强化学习控制的文献大多利用系统的输入输出数据,在系统模型未知的情况下学习整体控制器.但是,针对无人机系统存在标称控制器的情形,如何利用数据驱动的方法实现控制性能在线优化的研究较少.本文将Q学习方法^[9-10]和增量式控制结构相结合,提出一种增量式Q学习的无人机跟踪控制性能优化方法,在不修改预先设计的控制系统情况下,实现无人机控制性能优化.预先设计的控制系统为针对标称系统设计的控制器,由于模型误差和工作点变化等原因,其控制性能不一定最优,增量式控制器的作用为补偿控制器增益从而优化系统性能.

本文主要工作如下:1)提出一种无人机控制系统性能优化结构,不改变原有无人机控制器,利用Q

学习方法设计增量式控制器对标称控制器进行性能优化;2)分别针对无人机状态可知和状态不完全可知情况,设计基于状态反馈的增量式Q学习控制器与基于输出反馈的增量式Q学习控制器,并给出相应的策略迭代算法;3)在激励噪声满足持续激励条件下,证明了增量式Q学习方法进行参数估计的无偏性.

1 无人机建模与控制系统描述

1.1 无人机动力学模型

本文的仿真对象为F-16型固定翼飞行器.假设飞行器的质量恒视为常数,并且地面坐标系视为惯性系.在水平无侧滑和滚转的条件下,在地面坐标系中建立无人机的动力学方程,F-16的纵向平面模型描述^[11]为

$$\begin{cases} \dot{V} = \frac{T \cos \alpha - D}{m} + g(-\cos \alpha \sin \theta + \sin \alpha \cos \theta), \\ \dot{\alpha} = \frac{-T \sin \alpha - L}{mV} + q + \frac{g(\sin \alpha \sin \theta + \cos \alpha \cos \theta)}{V}, \\ \dot{\theta} = q, \\ \dot{q} = \frac{M}{I_y}. \end{cases} \quad (1)$$

其中: V 、 α 、 θ 和 q 分别表示无人机的速度、迎角、俯仰角和俯仰角速率; m 、 g 、 T 和 I_y 分别表示无人机的质量、重力加速度、发动机推力和转动惯量; M 、 D 和 L 分别为俯仰力矩、阻力和升力,表达式如下:

$$\begin{cases} M = \bar{q} S \bar{c} C_m, \\ D = \bar{q} S C_x, \\ L = \bar{q} S C_z, \end{cases} \quad (2)$$

\bar{q} 、 S 和 \bar{c} 分别表示动压、机翼面积和机翼平均气动弦长, C_x 、 C_z 和 C_m 分别表示阻力系数、升力系数和俯仰力矩系数.

本文使用F-16飞行器的低保真空气动力学模型.低保真模型忽略了飞机的前缘襟翼和速度制动器的作用.纵向模型的力和力矩系数计算如下:

$$\begin{cases} C_x = C_x(\alpha, \delta_e) + \frac{\bar{c}q}{2V} C_{xq}, \\ C_z = C_z(\alpha, \delta_e) + \frac{\bar{c}q}{2V} C_{zq}, \\ C_m = C_m(\alpha, \delta_e) + \frac{\bar{c}q}{2V} C_{mq} + C_z(x_{cgr} - x_{cg}), \end{cases} \quad (3)$$

其中 $C_x(\alpha, \delta_e)$ 、 $C_z(\alpha, \delta_e)$ 、 $C_m(\alpha, \delta_e)$ 、 C_{xq} 、 C_{zq} 和 C_{mq} 为气动参数并以插值表^[12]形式给出.这些空气动力学数据在 $-10^\circ \leq \alpha \leq 45^\circ$ 、 $0.1 \leq M_a \leq 0.6$ 飞行范围有效, M_a 为马赫数. x_{cgr} 、 x_{cg} 和 δ_e 分别为参考重心位

置、重心位置和升降舵偏转角. F-16飞行器的相关参数参见文献[12].

1.2 控制系统描述

将飞行器纵向模型(1)配平线性化后,考虑以下线性离散动力学系统:

$$\begin{aligned} x_{k+1} &= Ax_k + Bu_k, \\ y_k &= Cx_k. \end{aligned} \quad (4)$$

其中: $x_k = [V \ \alpha \ \theta \ q]^T \in \mathbf{R}^n$ 为系统状态, $u_k = \delta_e \in \mathbf{R}^m$ 为控制输入, $y_k = \alpha \in \mathbf{R}^p$ 为系统输出. 向量维数 $n = 4, m = 1, p = 1$. 参数矩阵 $A \in \mathbf{R}^{n \times n}, B \in \mathbf{R}^{n \times m}, C \in \mathbf{R}^{p \times n}$. 假设 (A, B) 可控, (A, C) 可观. 参考信号可由以下动力学系统生成:

$$\begin{aligned} x_{k+1}^r &= Fx_k^r, \\ y_k^r &= Cr_k^r. \end{aligned} \quad (5)$$

其中: $x_k^r \in \mathbf{R}^{n_1}$ 为参考系统状态, $y_k^r \in \mathbf{R}^p$ 为参考输出, 参数矩阵 $F \in \mathbf{R}^{n_1 \times n_1}, C_r \in \mathbf{R}^{p \times n_1}$. 结合式(4)和(5)可得增广系统^[13]如下:

$$X_{k+1} = TX_k + B_1u_k. \quad (6)$$

其中: $X_k = \begin{bmatrix} x_k \\ x_k^r \end{bmatrix} \in \mathbf{R}^{n+n_1}, T = \begin{bmatrix} A & 0 \\ 0 & F \end{bmatrix} \in \mathbf{R}^{(n+n_1) \times (n+n_1)}, B_1 = \begin{bmatrix} B \\ 0 \end{bmatrix} \in \mathbf{R}^{(n+n_1) \times m}$. 跟踪问题是将原系统与参考信号系统结合为增广系统,对增广系统进行研究并设计控制律.

为了优化系统性能,设计二次型形式的性能指标如下:

$$\begin{aligned} V(k) &= \sum_{i=k}^{\infty} \gamma^{i-k} [(y_i - y_i^r)^T Q (y_i - y_i^r) + u_i^T R u_i], \\ Q &\geq 0, R > 0. \end{aligned} \quad (7)$$

本文采用控制系统结构如图1所示,控制器包括预先设计的用来负责系统稳定的标称控制器和负责性能优化的增量式控制器. 增量式的性能优化不需改变标称控制器,而是产生一个增量式控制信号加入到控制输入端. 增量式控制器采用Q学习得到,该方法无需系统动力学模型,也无需预知预先设计的标称控制器增益大小. 在优化后,系统预先设计控制律与补偿控制律总和达到最优. 设计方法有如下两种:

1) 在系统全状态信息可获得的情况下,设计基于状态反馈的增量式Q学习控制器来完成系统对参考信号跟踪问题的性能优化.

2) 在系统全状态信息不能够获得的情况下,利用系统输入、输出数据信息重构状态变量,设计基于输

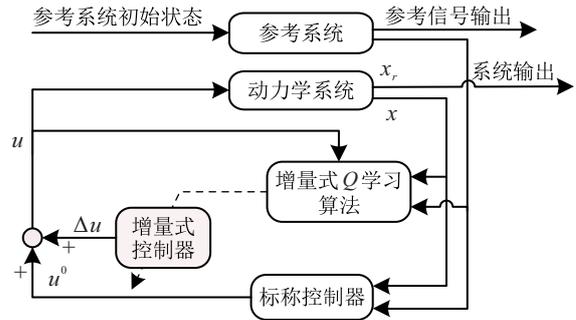


图1 总体控制框架

出反馈的增量式Q学习控制器优化跟踪性能.

两种补偿控制器设计方法均基于Q学习理论,无需知道系统动力学模型和预先设计的控制器增益大小. 在优化后,标称控制律与补偿控制律总和达到最优.

2 增量式Q学习状态反馈控制器设计

本节将Q学习方法与增量式控制相结合,介绍一种增量式Q学习方法. 通过Q学习去补偿现有控制器增益,从而对模型未知系统进行性能优化.

2.1 增量式Q学习与性能优化

考虑一种控制器结构,该控制器由预先设计的标称控制器和增量式控制器组成^[14],具体表示如下:

$$\begin{cases} u_k = u_k^0 + \Delta u_k, \\ u_k^0 = K_0 X_k, \\ \Delta u_k = \Delta K X_k. \end{cases} \quad (8)$$

其中: $K_0 \in \mathbf{R}^{m \times (n+n_1)}$ 表示预先设计的控制器增益,它的作用是使系统稳定,但在给定的性能指标下是非最优控制; $u_k^0 \in \mathbf{R}^m$ 表示标称控制器, $\Delta u_k \in \mathbf{R}^m$ 表示补偿控制器; $\Delta K \in \mathbf{R}^{m \times (n+n_1)}$ 表示需要设计的增量式控制器增益. 性能函数(7)等价于

$$V(X_k) = \sum_{i=k}^{\infty} \gamma^{i-k} (X_i^T Q_1 X_i + u_i^T R u_i). \quad (9)$$

其中: $Q_1 = \begin{bmatrix} C^T Q C & -C^T Q C_r \\ -C_r^T Q C & C_r^T Q C_r \end{bmatrix}$, γ 为折扣因子且 $0 < \gamma \leq 1$. $V(X_k)$ 的贝尔曼方程形式如下:

$$V(X_k) = X_k^T Q_1 X_k + u_k^T R u_k + \gamma V(X_{k+1}). \quad (10)$$

对于线性二次型问题, $V(X_k)$ 的二次型形式为

$$V(X_k) = X_k^T P X_k. \quad (11)$$

其中: P 为黎卡提方程的解,有

$$\begin{aligned} P &= Q_1 + \bar{T}^T P \bar{T} - \bar{T}^T P \bar{B}_1 (R + \bar{B}_1^T P \bar{B}_1)^{-1} \bar{B}_1^T P \bar{T}; \\ \bar{T} &= \gamma^{1/2} T, \bar{B}_1 = \gamma^{1/2} B_1. \end{aligned} \quad (12)$$

式(12)具有唯一解的充分必要条件是动力学系统 $(\gamma^{1/2} T, \gamma^{1/2} B_1)$ 是稳定的.

定义 Q 函数如下:

$$Q(X_k, u_k) = X_k^T Q_1 X_k + u_k^T R u_k + \gamma V(X_{k+1}). \quad (13)$$

Q 函数是依赖于状态和控制的函数.

将式(11)代入(13),有

$$\begin{aligned} Q(X_k, u_k) &= \\ & X_k^T Q_1 X_k + u_k^T R u_k + \gamma X_{k+1}^T P X_{k+1} = \\ & X_k^T Q_1 X_k + u_k^T R u_k + \gamma (T X_k + B_1 u_k)^T P \times \\ & (T X_k + B_1 u_k) = \\ & X_k^T Q_1 X_k + (K_0 X_k + \Delta u_k)^T R (K_0 X_k + \Delta u_k) + \\ & \gamma ((T + B_1 K_0) X_k + B_1 \Delta u_k)^T P \times \\ & ((T + B_1 K_0) X_k + B_1 \Delta u_k) := Q(X_k, \Delta u_k). \end{aligned}$$

整理得

$$Q(X_k, \Delta u_k) = \begin{bmatrix} X_k \\ \Delta u_k \end{bmatrix}^T H \begin{bmatrix} X_k \\ \Delta u_k \end{bmatrix}. \quad (14)$$

其中

$$T_F = T + B_1 K_0,$$

$$H =$$

$$\begin{bmatrix} Q_1 + K_0^T R K_0 + \gamma T_F^T P T_F & K_0^T R + \gamma T_F^T P B_1 \\ R K_0 + \gamma B_1^T P T_F & R + \gamma B_1^T P B_1 \end{bmatrix}.$$

为了使 Q 函数(13)最小化,执行 $\frac{\partial Q(X_k, \Delta u_k)}{\partial \Delta u_k} = 0$,可得最优控制律为

$$\begin{aligned} \Delta u_k &= \\ & -(R + \gamma B_1^T P B_1)^{-1} (R K_0 + \gamma B_1^T P T_F) X_k. \end{aligned} \quad (15)$$

式(15)等价于

$$\begin{aligned} \Delta u_k &= \\ & -(R + \gamma B_1^T P B_1)^{-1} (R K_0 + \\ & \gamma B_1^T P T + \gamma B_1^T P B_1 K_0) X_k = \\ & - \left(\frac{R}{\gamma} + B_1^T P B_1 \right)^{-1} B_1^T P T X_k - K_0 X_k. \end{aligned} \quad (16)$$

即总的控制输入为

$$\begin{aligned} u_k &= K_0 X_k + \Delta u_k = \\ & - \left(\frac{R}{\gamma} + B_1^T P B_1 \right)^{-1} B_1^T P T X_k. \end{aligned} \quad (17)$$

结合式(14),将核矩阵 H 表示为如下形式(即增量式控制器的无模型实现方式得以实现):

$$Q(X_k, \Delta u_k) = \begin{bmatrix} X_k \\ \Delta u_k \end{bmatrix}^T \begin{bmatrix} H_{XX} & H_{Xu} \\ H_{uX} & H_{uu} \end{bmatrix} \begin{bmatrix} X_k \\ \Delta u_k \end{bmatrix}.$$

其中: $H = H^T \in \mathbf{R}^{l \times l}$, $l = n + n_1 + m$; $H_{XX} = Q_1 + K_0^T R K_0 + \gamma T_F^T P T_F \in \mathbf{R}^{(n+n_1) \times (n+n_1)}$, H_{Xu}

$= K_0^T R + \gamma T_F^T P B_1 \in \mathbf{R}^{m \times (n+n_1)}$, $H_{uu} = R + \gamma B_1^T P B_1 \in \mathbf{R}^{m \times m}$. 则对应的无模型的增量式控制律为

$$\Delta u_k = -H_{uu}^{-1} H_{uX} X_k. \quad (18)$$

2.2 数据驱动实现与算法

对于式(8)所示控制器结构, Q 函数对应的贝尔曼方程为

$$\begin{aligned} z_k^T H z_k &= \\ & X_k^T Q_1 X_k + u_k^T R u_k + \gamma z_{k+1}^T H z_{k+1} = \\ & X_k^T (Q_1 + K_0^T R K_0) X_k + \Delta u_k^T R \Delta u_k + \\ & \Delta u_k^T R K_0 X_k + X_k^T K_0^T R \Delta u_k + \gamma z_{k+1}^T H z_{k+1} = \\ & z_k^T S_0 z_k + \gamma z_{k+1}^T H z_{k+1}. \end{aligned} \quad (19)$$

其中

$$\begin{aligned} z_k &= [X_k^T \quad \Delta u_k^T]^T \in \mathbf{R}^l, \\ S_0 &= \begin{bmatrix} Q_1 + \gamma K_0^T P K_0 & K_0^T R \\ R K_0 & R \end{bmatrix}. \end{aligned}$$

Q 函数矩阵 H 是未知的,需要通过 Q 学习算法学习得到,可通过参数化式(14)来分离 H ,即

$$Q(X_k, \Delta u_k) = z_k^T H z_k = \bar{H}^T \bar{z}_k. \quad (20)$$

其中: $\bar{H} \in \mathbf{R}^{l(l+1)/2}$, $\bar{z}_k \in \mathbf{R}^{l(l+1)/2}$,且满足

$$\begin{aligned} \bar{H} &\triangleq \\ & [h_{11}, 2h_{12}, \dots, 2h_{1l}, h_{22}, 2h_{23}, \dots, 2h_{2l}, \dots, h_{ll}]^T, \\ \bar{z}_k &\triangleq [z_1^2, z_1 z_2, \dots, z_1 z_l, z_2^2, z_2 z_3, \dots, z_2 z_l, \dots, z_l^2]^T. \end{aligned}$$

通过参数化式(19)后得到如下方程:

$$\bar{H}^T \bar{z}_k = z_k^T S_0 z_k + \gamma \bar{H}^T \bar{z}_{k+1}. \quad (21)$$

上述线性方程可以写成

$$\Phi^T \bar{H} = \Omega. \quad (22)$$

其中

$$\begin{aligned} \Phi &= [\bar{z}_k - \gamma \bar{z}_{k+1}, \bar{z}_{k+1} - \gamma \bar{z}_{k+2}, \dots, \bar{z}_{k+L-1} - \gamma \bar{z}_{k+L}], \\ \Omega &= [r_k, r_{k+1}, \dots, r_{k+L-1}]^T, r_k = z_k^T S_0 z_k, \\ \Phi &\in \mathbf{R}^{l(l+1)/2 \times L}, \Omega \in \mathbf{R}^L. \end{aligned}$$

式(22)的解为

$$\bar{H} = (\Phi \Phi^T)^{-1} \Phi \Omega. \quad (23)$$

在式(21)中, \bar{H} 是未知的向量,为求解线性方程,需要的样本数据需满足 $L \geq l(l+1)/2$. 下面提出策略迭代的强化学习算法来学习状态反馈增量式 Q 学习的补偿控制律.

算法1 状态反馈增量式 Q 学习策略迭代算法.

1) 初始化. 初始化 H^0 ,选取使系统稳定的标称

控制器 u^0 , 选择稳定的初始控制策略 Δu_0 , 对 $j(j = 1, 2, 3, \dots)$ 执行如下.

2) 策略评估. 收集足够步长的数据解Q函数贝尔曼方程

$$\begin{aligned} (\bar{H}^j)^T(\bar{z}_k - \gamma\bar{z}_{k+1}) = \\ (y_k - y_k^r)^T Q(y_k - y_k^r) + (u_k^j)^T R u_k^j. \end{aligned}$$

3) 策略更新. 在策略评估得到 \bar{H} 后构造 H , 计算改进后的控制策略

$$\Delta u_k^{j+1} = -(H_{uu}^j)^{-1} H_{uX}^j X_k.$$

4) 终止条件. 当达到收敛条件 $\|H^j - H^{j-1}\| \leq \varepsilon$ 时停止迭代, 其中 ε 为一个预先设计的正常数, 否则执行 $j = j + 1$ 继续迭代.

算法1中策略评估需要收集 $X_k, \Delta u_k, X_{k+1}$ 和 Δu_{k+1} 的数据求解贝尔曼方程(21), \bar{H} 的解如式(23)所示. 在策略更新步骤中, 通过最小化第 j 个策略的Q函数, 得到一个更好的策略 Δu_k^{j+1} . 在策略评估过程中, 需要加入激励噪声使得贝尔曼方程(21)有解, 激励噪声的作用阐述如下.

注1 通过在控制端加入满足持续激励条件的噪声, 使得向量 $\bar{z}_k - \gamma\bar{z}_{k+1}$ 线性无关, 从而使得线性方程(22)中的矩阵 Φ 满秩^[15]. 利用最小二乘法, 加入持续激励条件的噪声保证了式(23)解的唯一性.

定理1 激励噪声在求解Q函数贝尔曼方程时, 参数估计不产生偏差.

证明过程略.

为了保证式(19)的唯一解, 在控制输入端加入的激励噪声对Q函数的参数估计没有偏差作用, 学习到的增量式最优控制律不存在偏差. 算法1表明, 控制性能的优化可以通过增量的形式来实现, 而无需变动预先设计好的控制器增益.

注2 通过增量结构可使控制系统的控制器设计更加有效, 预先设计的标称控制器可能表现的是非最优性能, 它的主要作用是使系统稳定. 而增量式的控制器负责性能优化, 它可补偿地产生控制作用, 而无需更改原有控制系统.

3 增量式Q学习输出反馈控制器设计

3.1 增量式Q学习与性能优化

在算法1中, 需要系统的全状态信息来设计状态反馈控制器, 但状态信息在实际情况下可能获取不到. 本节将利用系统的输入、输出和参考信号序列的观测数据改进算法1, 改进的策略迭代算法利用观测数据对Q函数进行策略评估与与策略更新. 由文献[16]可知, 系统的状态可通过过去时刻测量到的输入、输出和参考信号数据来重构.

$$X_k = [M_u \ M_y \ M_r] \begin{bmatrix} \bar{u}_{k-1, k-N} \\ \bar{y}_{k-1, k-N} \\ x_{k-N}^r \end{bmatrix} = M\tau_k. \quad (24)$$

其中

$$\begin{aligned} \bar{u}_{k-1, k-N} &= [u_{k-1}^T \ u_{k-2}^T \ \dots \ u_{k-N}^T]^T, \\ \bar{y}_{k-1, k-N} &= [y_{k-1}^T \ y_{k-2}^T \ \dots \ y_{k-N}^T]^T. \end{aligned}$$

状态重构性表明, 在满足系统可观性条件 $\text{rank}(V_N) = n$ 下, 增广系统状态 X_k 可以用过去时刻系统 $[k - N, k - 1]$ 段的输入、输出以及参考信号数据来表示. 因此, 输出反馈下的控制器设计为

$$\begin{cases} u_k = u_k^0 + \Delta u_k, \\ u_k^0 = F_0 \tau_k, \\ \Delta u_k = \Delta F \tau_k, \end{cases} \quad (25)$$

其中 $F_0 = [F_u \ F_y \ F_r] \in \mathbf{R}^{m \times (mN + pN + n_1)}$. 利用该增广系统状态 X_k 的新表示形式和增量下输出的形式, 可以描述输入输出数据方案下的Q函数, 并在此基础上开发数据驱动的增量式Q学习算法. 为此, 记 $Z_k = [\bar{u}_{k-1, k-N}^T \ \bar{y}_{k-1, k-N}^T \ (x_{k-N}^r)^T \ \Delta u_k^T]^T \in \mathbf{R}^{l_1}$, $l_1 = mN + pN + n_1 + m$. 将 Z_k 代入式(13)中, Q函数表达式如下:

$$Q(Z_k) = Z_k^T H_1 Z_k. \quad (26)$$

$H_1 = H_1^T \in \mathbf{R}^{l_1 \times l_1}$ 可以分块为

$$H_1 = \begin{bmatrix} H_{\bar{u}\bar{u}} & H_{\bar{u}\bar{y}} & H_{\bar{u}r} & H_{\bar{u}u} \\ H_{\bar{y}\bar{u}} & H_{\bar{y}\bar{y}} & H_{\bar{y}r} & H_{\bar{y}u} \\ H_{r\bar{u}} & H_{r\bar{y}} & H_{rr} & H_{ru} \\ H_{u\bar{u}} & H_{u\bar{y}} & H_{ur} & H_{uu} \end{bmatrix}.$$

其中

$$\begin{aligned} H_{\bar{u}\bar{u}} &= M_u^T Q_1 M_u + F_u^T R F_u + \gamma G_u^T P G_u, \\ H_{\bar{u}\bar{y}} &= M_u^T Q_1 M_y + F_u^T R F_y + \gamma G_u^T P G_y, \\ H_{\bar{u}r} &= M_u^T Q_1 M_r + F_u^T R F_r + \gamma G_u^T P G_r, \\ H_{\bar{u}u} &= F_u^T R + \gamma G_u^T P B_1, \\ H_{\bar{y}\bar{y}} &= M_y^T Q_1 M_y + F_y^T R F_y + \gamma G_y^T P G_y, \\ H_{\bar{y}r} &= M_y^T Q_1 M_r + F_y^T R F_r + \gamma G_y^T P G_r, \\ H_{\bar{y}u} &= F_y^T R + \gamma G_y^T P B_1, \\ H_{rr} &= M_r^T Q_1 M_r + F_r^T R F_r + \gamma G_r^T P G_r, \\ H_{ru} &= F_r^T R + \gamma G_r^T P B_1, \\ H_{uu} &= R + \gamma B_1^T P B_1; \\ H_{\bar{u}\bar{u}} &\in \mathbf{R}^{mN \times mN}, \ H_{\bar{u}\bar{y}} \in \mathbf{R}^{mN \times pN}, \ H_{\bar{u}r} \in \mathbf{R}^{mN \times n_1}, \\ H_{\bar{u}u} &\in \mathbf{R}^{mN \times m}, \ H_{\bar{y}\bar{y}} \in \mathbf{R}^{pN \times pN}, \ H_{\bar{y}r} \in \mathbf{R}^{pN \times n_1}, \\ H_{\bar{y}u} &\in \mathbf{R}^{mN \times m}, \ H_{rr} \in \mathbf{R}^{n_1 \times n_1}, \ H_{ru} \in \mathbf{R}^{n_1 \times m}, \end{aligned}$$

$$H_{uu} \in \mathbf{R}^{m \times m}, G_u = TM_u + B_1F_u,$$

$$G_y = TM_y + B_1F_y, G_r = TM_r + B_1F_r.$$

式(26)将 Q 函数表示为输入、输出、参考信号和增量式输入的新形式,最优化控制策略通过 $\partial Q(Z_k)/\partial \Delta u_k = 0$ 获得.利用输入、输出和参考信号序列,可以得到最优的增量式控制器为

$$\Delta u_k = -H_{uu}^{-1} [H_{u\bar{u}} \ H_{u\bar{y}} \ H_{ur}] \begin{bmatrix} \bar{u}_{k-1,k-N} \\ \bar{y}_{k-1,k-N} \\ x_{k-N}^r \end{bmatrix} = \Delta F^* \begin{bmatrix} \bar{u}_{k-1,k-N} \\ \bar{y}_{k-1,k-N} \\ x_{k-N}^r \end{bmatrix}. \quad (27)$$

3.2 数据驱动实现与算法

在本节中,根据测量数据的 Q 函数-贝尔曼方程来评估控制策略

$$Z_k^T H_1 Z_k = (y_k - y_k^r)^T Q(y_k - y_k^r) + u_k^T R u_k + \gamma Z_{k+1}^T H_1 Z_{k+1}. \quad (28)$$

为了求解 H_1 来实现 Q 学习算法,现将 Q 函数(26)参数化为

$$Q(Z_k) = Z_k^T H_1 Z_k = \bar{H}_1^T \bar{Z}_k. \quad (29)$$

其中

$$\bar{H}_1 \triangleq [h_{11}, 2h_{12}, \dots, 2h_{1l_1}, h_{22}, 2h_{23}, \dots, 2h_{2l_1}, \dots, h_{l_1 l_1}]^T, \bar{Z} \triangleq [Z_1^2, Z_1 Z_2, \dots, Z_1 Z_{l_1}, Z_2^2, Z_2 Z_3, \dots, Z_2 Z_{l_1}, \dots, Z_{l_1}^2]^T.$$

上述线性方程可以写成

$$\Phi_1^T \bar{H}_1 = \Omega_1. \quad (30)$$

其中

$$\begin{aligned} \Phi_1 &= [\bar{Z}_k - \gamma \bar{Z}_{k+1}, \bar{Z}_{k+1} - \gamma \bar{Z}_{k+2}, \dots, \\ &\quad \bar{Z}_{k+L-1} - \gamma \bar{Z}_{k+L}], \\ \Omega_1 &= [r_k, r_{k+1}, \dots, r_{k+L-1}]^T, \\ r_k &= (y_k - y_k^r)^T Q(y_k - y_k^r) + u_k^T R u_k, \\ \Phi_1 &\in \mathbf{R}^{l_1(l_1+1)/2 \times L}, \Omega_1 \in \mathbf{R}^L. \end{aligned}$$

式(30)的解为

$$\bar{H}_1 = (\Phi_1 \Phi_1^T)^{-1} \Phi_1 \Omega_1. \quad (31)$$

在式(29)中, \bar{H}_1 是未知的向量,为了解线性方程,需要的样本数据需满足 $L \geq l_1(l_1 + 1)/2$.下面提出策略迭代的强化学习算法来学习输出反馈增量式 Q 学习的补偿控制律.

算法2 输出反馈增量式 Q 学习策略迭代算法.

1) 初始化. 初始化 \bar{H}_1^0 ,选取使系统稳定的标称控制器 u^0 ,选择稳定的初始控制策略 Δu_0 ,对 $j(j = 1, 2, 3, \dots)$ 执行如下操作.

2) 策略评估. 收集足够步长的数据解 Q 函数贝尔曼方程

$$\begin{aligned} &(\bar{H}_1^j)^T (\bar{Z}_k - \gamma \bar{Z}_{k+1}) = \\ &(y_k - y_k^r)^T Q(y_k - y_k^r) + (u_k^j)^T R u_k^j. \end{aligned}$$

3) 策略更新. 在策略评估得到 \bar{H}_1 后构造 H_1 ,计算改进后的控制策略

$$\begin{aligned} \Delta u_k^{j+1} &= -(H_{uu}^j)^{-1} (H_{u\bar{u}}^j \bar{u}_{k-1,k-N} + \\ &H_{u\bar{y}}^j \bar{y}_{k-1,k-N} + H_{ur}^j x_{k-N}^r). \end{aligned}$$

4) 终止条件. 当达到收敛条件 $\|\bar{H}_1^j - \bar{H}_1^{j-1}\| \leq \varepsilon$ 时停止迭代,其中 ε 为一个预先设计的正常数,否则执行 $j = j + 1$ 继续迭代.

注3 算法1和算法2所示的增量式控制律是一种无模型的方法,只需获得辨识核矩阵的样本数据,便可在数据驱动的环境下自适应地学习到该增量控制律.由于系统的标称控制律使系统稳定,在算法中使用策略迭代算法来求解 H 和 H_1 .

4 仿真分析

本文使用F-16飞行器模型进行算法验证.选择使用低保真空气动力模型,在高度为5000ft和速度为300ft/s的平稳飞行状态下配平.离散时间取为0.5s,离散化后F-16纵向模型的线性离散系统描述如下:

$$\begin{aligned} x_{k+1} &= \begin{bmatrix} 0.9854 & 0.6528 & -15.9677 & -4.5261 \\ -0.0003 & 0.7891 & 0.0026 & 0.3305 \\ 0 & 0.0332 & 1 & 0.4153 \\ 0 & 0.1191 & 0.0001 & 0.6912 \end{bmatrix} x_k + \\ &\begin{bmatrix} -0.0483 \\ 0.0057 \\ 0.0063 \\ 0.0237 \end{bmatrix} u_k, \end{aligned}$$

$$y_k = [0 \ 57.3 \ 0 \ 0] x_k. \quad (32)$$

其中: $x_k = [V \ \alpha \ \theta \ q]^T, u = \delta_e, \delta_e$ 的偏转范围为 $\pm 25.0^\circ$.

参考信息由以下动力学系统^[17]生成:

$$x_{k+1}^r = \begin{bmatrix} 0.9954 & -0.0954 \\ 0.0954 & 0.9954 \end{bmatrix} x_k^r, y_k^r = [5 \ 0] x_k^r. \quad (33)$$

4.1 增量式状态反馈

各个状态的初始值设置为 $V = 10 \text{ ft/s}$, $\alpha = 10^\circ$, $\theta = -10^\circ$, $q = 10^\circ/\text{s}$, 参考轨迹系统初始状态设置为 $[1, -1]^T$, 代价函数中的权重矩阵 $R = 1$, $Q = 20$, 折扣因子 $\gamma = 0.3$, $H^0 = 0.01I$. 由于系统(32)是不稳定的, 需要预先设计标称控制器使其稳定, 标称控制器增益设计如下:

$$K_0 = [-0.0061 \ 6 \ 0.00794 \ 3 \ -1 \ 0.1784]. \quad (34)$$

在使用算法1学习最优增量式控制器时, 需要辨识的 H 维度为 7×7 , 则辨识一次需要满足采样数据长度 $L \geq 28$, 这里 L 取28, 通过算法1可得

$$\begin{aligned} H_{uX} &= [0 \ -0.0105 \ -0.0001 \ -0.0053], \\ H_{ur} &= [0.0011 \ -0.0001], \quad H_{uu} = 0.0002. \end{aligned} \quad (35)$$

将式(35)代入(18), 可得最优增量式控制律为

$$\begin{aligned} \Delta K &= [-0.0203 \ 47.4610 \ 0.2546 \ \rightarrow \\ &\leftarrow 23.9667 \ -5.1579 \ 0.6034]. \end{aligned} \quad (36)$$

将标称控制律(34)与学习到的增量式控制律(36)相加可得总控制律为

$$\begin{aligned} K_{\text{sum}} &= [-0.0264 \ 53.4610 \ 0.3343 \ \rightarrow \\ &\leftarrow 26.9667 \ -6.1579 \ 0.7818]. \end{aligned} \quad (37)$$

通过非增量式Q学习得到的控制律^[18]为

$$\begin{aligned} K_{\text{sum}}^* &= [-0.0264 \ 53.4692 \ 0.3323 \ \rightarrow \\ &\leftarrow 26.9678 \ -6.1588 \ 0.7752]. \end{aligned} \quad (38)$$

比较式(37)与(38)可知, 增量式Q学习能够在系统预先有一个标称控制器的情况下, 自适应地补偿标称控制律, 使得总的控制律达到最优.

状态反馈仿真如图2~图6所示. 图2表示迎角的变化曲线, 可以看出, 在标称控制器下, 迎角虽然不发散, 但对参考信号跟踪效果不理想, 在算法1的控制律补偿下, 迎角在第40步能够跟踪上参考信号且跟踪效果显著优于未优化时的效果. 图3为状态反馈下的升降舵偏转角曲线, 其范围为 $[-4^\circ, 10^\circ]$. 无人机配平点的升降舵偏转角为 -4.2383° , 因此, 总的升降

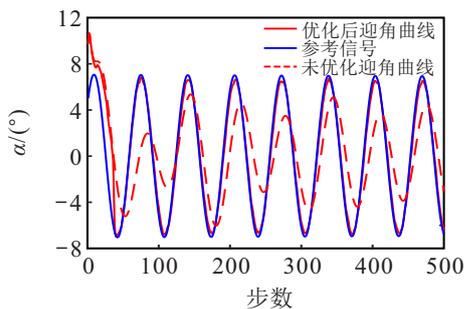


图2 状态反馈下迎角变化曲线

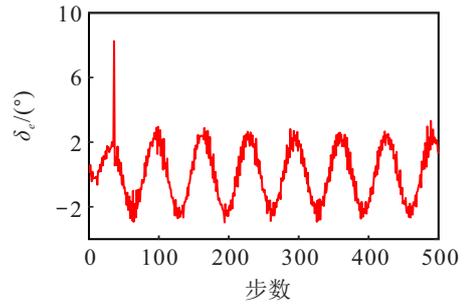


图3 状态反馈下升降舵偏转角

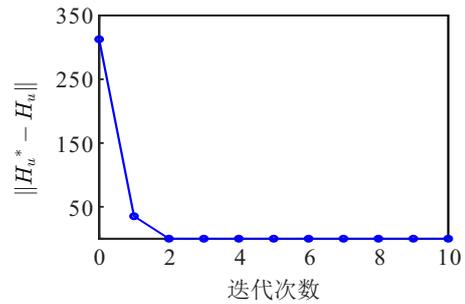


图4 状态反馈下H的收敛性

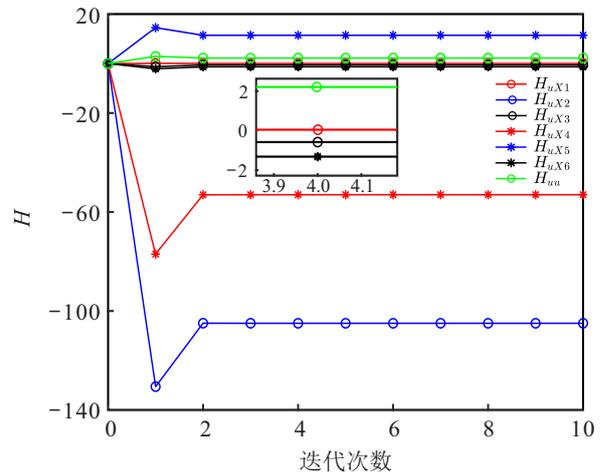


图5 状态反馈下H元素的收敛性

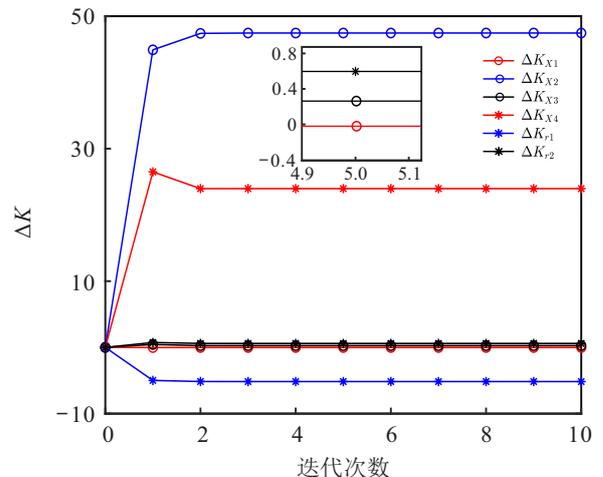


图6 状态反馈下控制律增益收敛性

舵偏转角在偏转角约束范围 $\pm 25.0^\circ$ 内. 图4和图5表示 H 中与增量式控制律有关的各元素收敛性, 在第2次迭代中即可学习到最优的增量式控制律. 图6为状

态反馈下增量式控制器增益收敛性.

4.2 增量式输出反馈

各个状态的初始值设置为 $V = 10 \text{ ft/s}$, $\alpha = 10^\circ$, $\theta = -10^\circ$, $q = 10^\circ/\text{s}$, 参考轨迹系统初始状态设置为 $[1, -1]^T$, 代价函数中的权重矩阵 $R = 1$, $Q = 7$, 折扣因子 $\gamma = 0.5$, $H_1^0 = 0.01I$. 预先设计如下标称控制器使系统(32)稳定:

$$\begin{aligned} F_0 &= [F_u \ F_y \ F_r], \\ F_u &= [-0.2174 \ 0.5217 \ 0.0609 \ -0.2522], \\ F_y &= [2.4348 \ -4.2174 \ 2.6522 \ -0.5739], \\ F_r &= [-1.3217 \ 0.7043]. \end{aligned} \quad (39)$$

在使用算法2学习最优增量式控制器时, 需要辨识的 H_1 维度为 11×11 , 则辨识一次需要满足采样数据长度 $L \geq 66$, 这里 L 取 70, 通过算法2可得

$$\begin{aligned} H_{u\bar{u}} &= [1.2444 \ -4.5613 \ -0.1932 \ 2.2322], \\ H_{u\bar{y}} &= [-18.2451 \ 34.6634 \ -22.7891 \ 5.0791], \\ H_{u\bar{r}} &= [6.3711 \ -3.8987], \ H_{ur} = 2.1107. \end{aligned} \quad (40)$$

将式(40)代入(27), 可得最优增量式控制律为

$$\begin{aligned} \Delta F_u &= [-0.5896 \ 2.1610 \ 0.0915 \ -1.0575], \\ \Delta F_y &= [8.6439 \ -16.3182 \ 10.7968 \ -2.4063], \\ \Delta F_r &= [-3.0184 \ 1.8471]. \end{aligned} \quad (41)$$

将标称控制律(40)与学习到的增量式控制律(41)相加, 可得总控制律为

$$\begin{aligned} F_{u_sum} &= [-0.8070 \ 2.6827 \ 0.1524 \ -1.3097], \\ F_{y_sum} &= [11.0787 \ -20.5356 \ 13.4490 \ -2.9802], \\ F_{r_sum} &= [-4.3401 \ 2.5514]. \end{aligned} \quad (42)$$

通过非增量式 Q 学习而学习到的控制律^[18]为

$$\begin{aligned} F_u^* &= [-0.8069 \ 2.6869 \ 0.1511 \ -1.3126], \\ F_y^* &= [11.0914 \ -20.5673 \ 13.4743 \ -2.9867], \\ F_r^* &= [-4.3404 \ 2.5512]. \end{aligned} \quad (43)$$

比较式(42)与(43)可知, 输出反馈增量式 Q 学习能够在系统预先有一个标称控制器的情况下, 利用输入输出数据设计增量式控制律, 自适应地补偿标称控制律, 使得总的控制律达到最优.

输出反馈仿真如图7~图11所示. 图7表示迎角的变化曲线, 可以看出, 标称控制器对参考信号跟踪效果不理想, 在算法2的控制律补偿下, 迎角在第100步能够跟踪上参考信号且跟踪效果显著优于未优化时的效果, 输出反馈迎角跟踪上参考信号的时间长于状态反馈, 是由于辨识一次 H 所需的样本数据步

长不同, 算法2所需数据步长明显多于算法1. 图8为输出反馈下的升降舵偏转角, 其范围为 $[-8^\circ, 6^\circ]$. 根据无人机配平点的升降舵偏转角为 -4.2383° , 因此, 总的升降舵偏转角在偏转角约束范围 $\pm 25.0^\circ$ 内. 图9和图10表示 H_1 中与增量式控制律有关的各元素收敛性, 在第2次迭代中即可学习到最优的增量式控制律. 图11表示增量式控制器增益收敛性.

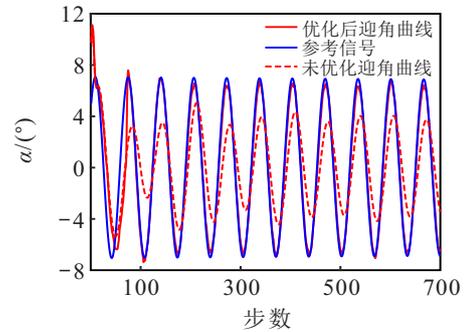


图7 输出反馈下迎角变化曲线

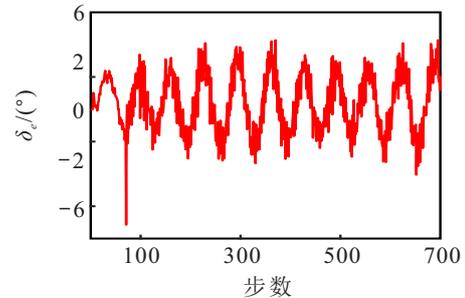


图8 输出反馈下升降舵偏转角

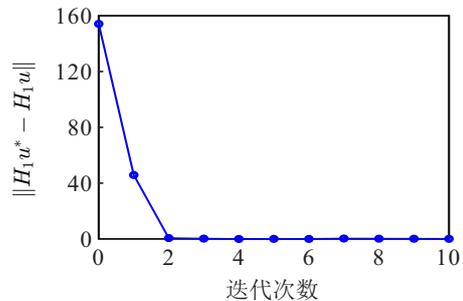


图9 输出反馈 H_1 的收敛性

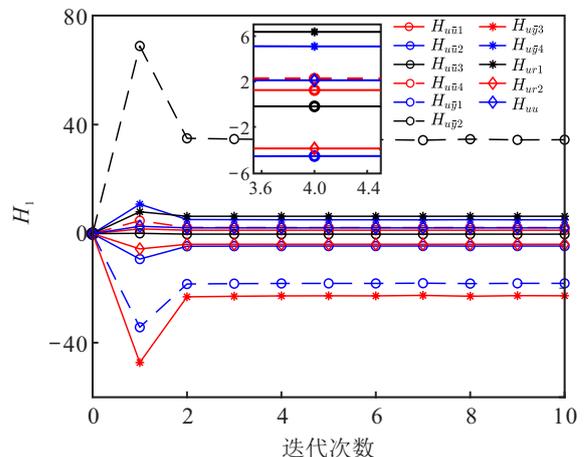


图10 输出反馈 H_1 元素的收敛性

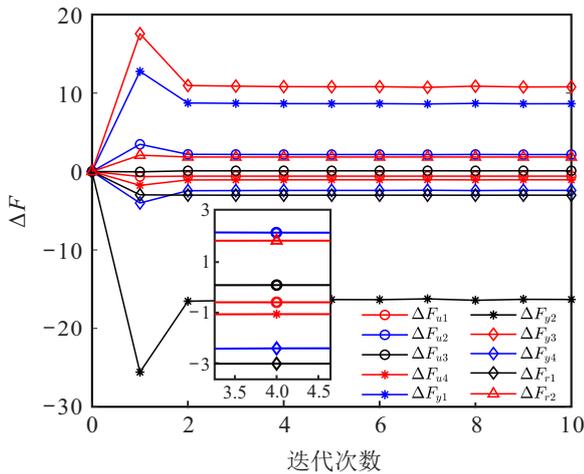


图 11 输出反馈控制律增益收敛性

注4 当激励噪声满足持续激励条件时,数据驱动的策略迭代算法的收敛性和利用李雅普诺夫方程迭代求解黎卡提方程的收敛性一致^[14]. 由于采用迭代算法的控制器增益满足二次收敛特性^[19],本文采用的Q学习算法也具有二次收敛特性,其收敛速度较快. 在二次收敛特性下,Q学习算法的参数设置主要包括权重矩阵 Q, R ,折扣因子 γ ,初始矩阵 H^0 和每步策略评估的数据长度 N . 当数据长度 N 满足矩阵满秩条件时,增大每步的数据长度不影响收敛速度. 因此,Q学习的收敛速度主要与参数 Q, R, γ, H^0 的设置有关.

4.3 数字PID控制器性能对比

对于本文的迎角跟踪,PID控制系统包含迎角偏差控制通道,迎角偏差控制通道控制器为

$$e_k = y_k - y_k^r$$

$$u_k = K_p e_k + K_i \sum_{j=0}^k e_j + K_d (e_k - e_{k-1}). \quad (44)$$

通过调试,PID控制参数为 $K_p = 0.5, K_i = 0.34, K_d = 0.25$. 将PID控制与本文方法进行对比,如图12和图13所示. 图12是增量式Q学习状态反馈优化控制与PID控制性能对比图,图13是增量式Q学习输出反馈优化控制与PID控制性能对比图,三者都能实现迎角对参考信号的跟踪,但本文提出的增量式Q学习算法比PID控制跟踪误差较小,且在参考信号跟踪初

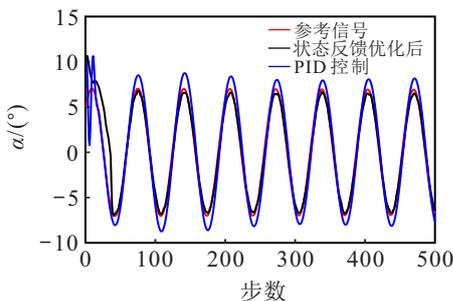


图 12 状态反馈与PID性能对比

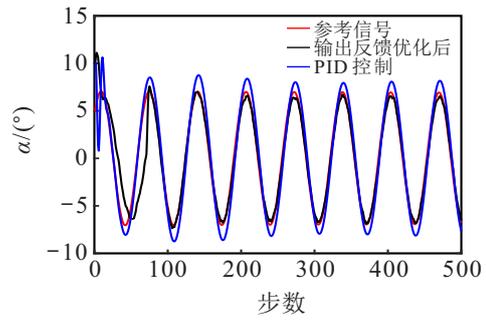


图 13 输出反馈与PID性能对比

始段比PID控制的响应更为平稳. 更重要的是,当无人机的工况点发生变化时,PID控制不能自适应地调整,其控制性能也会随着下降,而增量式Q学习控制可以根据工况点的变化自适应地学习一个新的控制器,实现无人机控制性能的实时优化. 因此,增量式Q学习控制器相比PID控制具有更优良的性能.

5 结论

本文针对固定翼无人机纵向控制的高性能需求,提出一种控制系统性能优化结构. 该控制结构在不修改预先设计的控制系统的情况下,通过优化一个增量式的控制系统增益来优化性能. 增量式控制器的设计分为两种:1)在无人机状态信息可获得情况下,设计基于状态反馈的增量式Q学习的补偿控制律;2)当无人机状态信息不能完全获得时,利用输入、输出和参考信号数据设计基于输出反馈的增量式Q学习的补偿控制律. 两种设计方法均为数据驱动下的无模型方法,且激励噪声对Q函数贝尔曼方程的参数估计没有偏差. 通过对选用的F-16纵向模型的仿真可知,增量式Q学习算法会自适应地补偿标称控制增益使总的控制输入达到最优. 与标称控制性能相比,增量式Q学习控制器的跟踪控制性能得到了较好的优化.

参考文献(References)

- [1] 王美仙, 李明, 张子军. 飞行器控制律设计方法发展综述[J]. 飞行力学, 2007, 25(2): 1-4.
(Wang M X, Li M, Zhang Z J. Developing status of control law design methods for flight[J]. Flight Dynamics, 2007, 25(2): 1-4.)
- [2] 江琼, 陈怀民, 吴佳楠. H_∞ 鲁棒控制与PID控制相结合的无人机飞行控制研究[J]. 宇航学报, 2006, 27(2): 192-195.
(Jiang Q, Chen H M, Wu J N. Research on UAV flight control based on PID control and H_∞ robust control[J]. Journal of Astronautics, 2006, 27(2): 192-195.)
- [3] 刘燕斌, 陆宇平. 基于反步法的高超音速飞机纵向逆飞行控制[J]. 控制与决策, 2007, 22(3): 313-317.

- (Liu Y B, Lu Y P. Longitudinal inversion flight control based on backstepping for hypersonic vehicle[J]. Control and Decision, 2007, 22(3): 313-317.)
- [4] Hwangbo J, Sa I, Siegwart R, et al. Control of a quadrotor with reinforcement learning[J]. IEEE Robotics and Automation Letters, 2017, 2(4): 2096-2103.
- [5] Sugimoto T, Gouko M. Acquisition of hovering by actual UAV using reinforcement learning[C]. The 3rd International Conference on Information Science and Control Engineering. Beijing: IEEE, 2016: 148-152.
- [6] Wang Y D, Sun J, He H B, et al. Deterministic policy gradient with integral compensator for robust quadrotor control[J]. IEEE Transactions on Systems, Man, and Cybernetics: Systems, 2020, 50(10): 3713-3725.
- [7] Zhou Y, Kampen E J V, Chu Q P. Nonlinear adaptive flight control using incremental approximate dynamic programming and output feedback[J]. Journal of Guidance, Control, and Dynamics, 2016, 40(2): 493-496.
- [8] de Alvear Cárdenas J I, Sun B, Van Kampen E J. Intelligent adaptive control using LADP and IADP applied to F-16 aircraft with imperfect measurements[C]. AIAA Scitech 2021 Forum. Reston: AIAA, 2021: 1119.
- [9] Kiumarsi B, Lewis F L, Modares H, et al. Reinforcement Q -learning for optimal tracking control of linear discrete-time systems with unknown dynamics[J]. Automatica, 2014, 50(4): 1167-1175.
- [10] Kiumarsi B, Lewis F L, Naghibi-Sistani M B, et al. Optimal tracking control of unknown discrete-time linear systems using input-output measured data[J]. IEEE Transactions on Cybernetics, 2015, 45(12): 2770-2779.
- [11] Nguyen L T. Simulator study of stall/post-stall characteristics of a fighter airplane with longitudinal static stability[M]. Washington, DC: National Aeronautics and Space Administration, Scientific and Technical Information Branch, 1979: 36-43.
- [12] Stevens B L, Lewis F L, Johnson E N. Aircraft control and simulation: Dynamics, controls design, and autonomous systems[M]. The 3rd edition. Hoboken: Wiley, 2015: 584-592.
- [13] Peng Y J, Chen Q, Sun W J. Reinforcement Q -learning algorithm for H_∞ tracking control of unknown discrete-time linear systems[J]. IEEE Transactions on Systems, Man, and Cybernetics: Systems, 2020, 50(11): 4109-4122.
- [14] Xu Y S, Zhao Z G, Yin S. Performance optimization and fault-tolerance of highly dynamic systems via Q -learning with an incrementally attached controller gain system[J]. IEEE Transactions on Neural Networks and Learning Systems, 2022 (99): 1-11.
- [15] Lewis F L, Vamvoudakis K G. Reinforcement learning for partially observable dynamic processes: Adaptive dynamic programming using measured output data[J]. IEEE Transactions on Systems, Man, and Cybernetics — Part B: Cybernetics, 2011, 41(1): 14-25.
- [16] Rizvi S A A, Lin Z L. Output feedback Q -learning for discrete-time linear zero-sum games with application to the H_∞ control[J]. Automatica, 2018, 95: 213-221.
- [17] Peng Y J, Meng Q Q, Sun W J. Adaptive output-feedback quadratic tracking control of continuous-time systems via value iteration with its application[J]. IET Control Theory & Applications, 2020, 14(20): 3621-3631.
- [18] Sun W J, Zhao G Y, Peng Y J. Adaptive optimal output feedback tracking control for unknown discrete-time linear systems using a combined reinforcement Q -learning and internal model method[J]. IET Control Theory & Applications, 2019, 13(18): 3075-3086.
- [19] Hwer G. An iterative technique for the computation of the steady state gains for the discrete optimal regulator[J]. IEEE Transactions on Automatic Control, 1971, 16(4): 382-384.

作者简介

赵振根(1989—),男,讲师,博士,从事信息物理系统安全、无人机智能控制等研究, E-mail: zhaozhengen@nuaa.edu.cn;

程磊(1997—),男,硕士生,从事无人机智能控制的研究, E-mail: chenglei@nuaa.edu.cn.