

文章编号: 1001-0920(2013)06-0930-05

操作条件反射学习自动机及其在机器人平衡控制中的应用

郜园园, 阮晓钢, 宋洪军

(北京工业大学 电子信息与控制工程学院, 北京 100124)

摘要: 针对两轮机器人的平衡控制问题, 在学习自动机理论的框架中, 提出一种基于操作条件反射学习自动机的仿生学习模型. 该模型引入认知学习单元和取向单元, 分别用来实现操作行为学习和指导系统进化的方向. 模拟两轮自平衡机器人的平衡控制仿真实验表明, 该学习模型具有可行性和有效性, 能使机器人自主学习平衡控制技能, 并使其具有高度的自适应能力.

关键词: 操作条件反射; 学习自动机; 仿生; 机器人平衡控制

中图分类号: TP273

文献标志码: A

Operant conditioning learning automatic and its application on robot balance control

GAO Yuan-yuan, RUAN Xiao-gang, SONG Hong-jun

(College of Electronic Information and Control Engineering, Beijing University of Technology, Beijing 100124, China. Correspondent: GAO Yuan-yuan, E-mail: yuangao84@163.com)

Abstract: A biomimetic learning model is proposed based on the operant conditioning learning automatic(OCLA) within the structure of learning automatic theory for the balance control of a two-wheeled robot. A cognitive learning unit and a tropism unit are introduced in this model, and they are used to evaluate the operant behavior learning and to direct the evolution of the system respectively. Finally, the simulation experiments on the two-wheeled robot show the feasibility and effectiveness of the proposed algorithm, and the robot learns the ability of balance control and has high self-adaptive ability.

Key words: operant conditioning; learning automatic; biomimetic; robot balance control

0 引言

两轮自平衡机器人是典型的欠驱动系统, 近年来, 针对两轮机器人的运动平衡控制问题已提出了较多控制方法, 如PID控制、LQR控制、模糊自适应控制、鲁棒控制和强化学习等. 但是, 这些方法均以传统的控制理论和控制技术为主, 导致动态性能和静态性能都较差.

操作条件反射(OC)^[1]是一种重要的条件反射理论, 视为生物系统最基本的学习形式, 这是因为人或动物的平衡控制技能在较大程度上是基于这种学习机制自组织地渐近形成、发展和完善的. 国内外许多学者对于OC的仿生学习模型进行了相关研究, 期望这种模型能够复制动物学习操作或控制的实验. 也有学者以概率自动机为平台模拟操作条件反射机制, 设计了相应的仿生系统, 并成功实现了倒立摆和机器人

的平衡控制^[2-3]. 但是, 这些计算理论和计算模型没有给出具体的数学计算模型, 不具备泛化能力, 应用受到限制^[4-7], 所以研究具有更加普适性的操作条件反射仿生学习模型具有重要的研究价值.

学习自动机(LA)可以作为机器人特别是认知发育机器人的数学抽象和形式化的工具, 近年来在机器人学中的应用研究逐渐增多. LA可以描述机器人及其环境^[8], 用于机器人环境勘测、绘制环境地图^[9], 帮助机器人进行行为选择^[10]和最优控制^[11]等. 针对两轮机器人的平衡控制问题, 以学习自动机为框架建立操作条件反射机制的数学模型是一个重要的研究思路, 也是本文的主要研究内容.

基于上述现状, 本文将学习自动机与操作条件反射的思想相结合, 设计一种仿生的学习模型, 并将其应用于两轮机器人的平衡控制问题. 认知学习单元主

收稿日期: 2012-02-01; 修回日期: 2012-05-02.

基金项目: 国家自然科学基金项目(61075110); 国家863计划项目(2007AA04Z226); 北京市自然科学基金项目(4102011); 北京市教委重点项目(KZ201210005001).

作者简介: 郜园园(1984—), 女, 博士生, 从事智能控制、机器学习的研究; 阮晓钢(1958—), 男, 教授, 博士生导师, 从事机器人、神经网络等研究.

$p_{ijk}(t)$ 的分布, 随机选择 t 时刻的操作 $o_k(t)$.

Step 3: 实施操作行为. 实施选取的操作行为 $o_k(t)$, 并依据 $f_S : S(t) \times A(t) \times O(t) \rightarrow S(t+1)$ 状态转移方程观测 $t+1$ 时刻 OCLA 的状态.

Step 4: 取向值增量计算. 依据取向单元中取向函数 ψ , 分别计算状态 $s(t)$ 和 $s(t+1)$ 的取向值, 得到取向值增量函数 $\psi(\cdot)$.

Step 5: 操作条件反射. 由认知学习单元 δ 按照式 (1) 对随机“条件-操作”规则 r_{ijk} 激发概率 $p_{ijk}(t)$ 进行调节.

Step 6: 对外输出. 由系统的输出方程 $f_Z : S(t) \times A(t) \times O(t) \rightarrow Z(t+1)$ 对外输出 $Z(t+1)$.

Step 7: 条件停止. 重复 Step 2~Step 6, 直至达到迭代学习次数 T_f 或 $p_{abc}(t+1) > p_\varepsilon$.

2 算法收敛性分析

引理 1 设 $OCLA = \langle A, S, O, Z, R, f, \psi, \delta \rangle$ 是一个操作条件反射学习自动机, 其状态转移过程 $f : S \times A \times O \rightarrow S, S \times A \times O \rightarrow Z$ 是确定的, 设 $p \neq 0$ 和 O_i^+ 是系统处于状态 s_i 下表现为正取向性 ($\psi_{ijk} > 0$) 的操作行为集合, 且 $O_i^+(t=0) \neq \emptyset, O_i^+ \subset O$, 则有

$$p_{ijk} = p(O_i^+ | s_i, a_j) = 1, t \rightarrow \infty. \quad (2)$$

即当 $t \rightarrow \infty$ 时, OCLA 依概率 1 选取正取向状态转移操作.

证明 令 l 为 s_i 出现次数的序号, $p_{ijk}(l)$ 为状态 s_i 第 l 次出现时集合 O_i^+ 中随机操作被选中的概率, 初始概率为 $p_{ijk}(0)$.

因为 $O_i^+(t=0) \neq \emptyset$, 所以 $p_{ijk}(0) \neq 0$, 又有 $\Delta p_{ijk}(l) \geq 0$, 故 $\forall l=0, 1, \dots, p_{ijk}(l) \neq 0$. 因为 $p_{ijk}(l) \neq 0$, 当 $t \rightarrow \infty$ 时, OCLA 状态 s_i 出现的频次趋于无穷, 同时 $p_{ijk}(l)$ 使得 O_i^+ 中操作行为被选中的频次也趋于无穷.

假定 OCLA 在第 l 次处于状态 s_i 和输入 a_j 的条件下选择操作 $o_k \in O$, 若 $o_k \in O_i^+$, 则 $\psi_{ijk} > 0, \xi(\psi_{ijk}) > 0$, 所以有

$$\begin{aligned} \Delta p_{ijk}(l) &= p_{ijk}(l+1) - p_{ijk}(l) = \\ &\xi(\psi_{ijk}) \sum_{m \neq k} p_{ijm}(l) \geq 0. \end{aligned} \quad (3)$$

$\Delta p_{ijk}(l) \geq 0$ 表示当 $t \rightarrow \infty$ 时 O_i^+ 被选中的频次趋于无穷. 当 $t \rightarrow \infty$ 时, $\Delta p_{ijk}(l) > 0$ 的情形可发生任意多次, 因为 $p_{ijk}(l)$ 有上界且为 1, $p_{ijk}(l)$ 增加到 1 为止. \square

定义 1 设 $OCLA = \langle A, S, O, Z, R, f, \psi, \delta \rangle$ 是一个操作条件反射自动机, 其操作熵 H 可定义为

$$H = H(O|S, A) =$$

$$\begin{aligned} \sum_{i=0}^{N_S} p_i H_i &= \sum_{i=0}^{N_S} \sum_{j=0}^{N_A} p(s_i, a_j) H_i(O|s_i, a_j) = \\ &= \sum_{i=0}^{N_S} \sum_{j=0}^{N_A} p(s_i, a_j) \sum_{k=1}^{N_O} p(o_k | s_i, a_j) \log_2 p(o_k | s_i, a_j), \end{aligned} \quad (4)$$

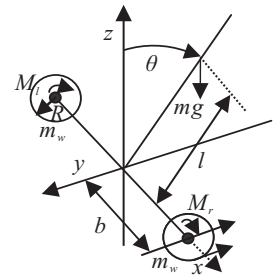
其中 H_i 为 OCLA 处于状态 s_i 条件下的操作熵.

3 基于 OCLA 的两轮机器人平衡控制

两轮自平衡机器人是本质上多变量、强耦合和非线性的复杂动态系统, 其核心问题是运动平衡控制. 本文以北工大人工智能与机器人研究所研制的两轮机器人为研究对象, 机器人系统见图 2, 主要物理参数见表 1.



(a) 实物图



(b) 受力分析

图 2 两轮机器人系统

表 1 两轮机器人的主要物理参数

物理量	符号	数值
机身的质量	m	10 kg
机身高度	H	0.65 m
质心距	l	0.35 m
车体绕竖直轴的转动惯量	I_m	0.589 3 kg·m ²
左右轮子的质量	m_w	1 kg
轮子半径	R	0.15 m
轮轴长度	$2b$	0.44 m
左右轮的转动惯量	I_w	0.011 3 kg·m ²
重力加速度	g	9.8 m/s ²
直流电机电阻	R_a	0.317 18 Ω
直流电机反电动势系数	K_e	0.030 6 V·s/rad
直流电机电磁转矩系	K_m	0.030 2 N·m/A
减速比	N	28

分别对机器人车轮、本体进行受力分析, 以牛顿经典力学为基础建立两轮机器人的数学模型. 通过对电机建模, 得到左右轮电机输出转矩 M_l 和 M_r 与电机控制电压 u_l 和 u_r 之间的关系, 分别为

$$M_l = -\frac{0.105 K_m^2 N^2}{R_a R} \dot{x} + \frac{N K_m}{R_a} u_l, \quad (5)$$

$$M_r = -\frac{0.105 K_m^2 N^2}{R_a R} \dot{x} + \frac{N K_m}{R_a} u_r. \quad (6)$$

通过式 (6) 和 (7) 将左右轮转矩控制转化为电压控制. 设 $X = [x \ \dot{x} \ \theta \ \dot{\theta}]^T$, 线性化后经过整理得到系统的状态空间表达式

$$\dot{X} = \begin{bmatrix} \dot{x} \\ \dot{\ddot{x}} \\ \dot{\theta} \\ \dot{\ddot{\theta}} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & -b_1 b_3 / a_2 & a_1 / a_2 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & b_2 b_3 / a_2 & a_3 / a_2 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ b_1 b_4 / a_2 & b_1 b_4 / a_2 \\ 0 & 0 \\ -b_2 b_4 / a_2 & -b_2 b_4 / a_2 \end{bmatrix} \begin{bmatrix} u_l \\ u_r \end{bmatrix}$$

其中

$$\begin{aligned} a_1 &= -2m^2 l^2 R^2 g, \\ a_2 &= (2ml^2 + I_m)(mR^2 + 2m_w R^2 + 2I_w) - 2m^2 l^2 R^2, \\ a_3 &= mgl(mR^2 + 2m_w R^2 + 2I_w), \\ b_1 &= 2ml^2 R + I_m R, \quad b_2 = mlR, \\ b_3 &= \frac{0.21K_m^2 N^2}{R_a R}, \quad b_4 = \frac{NK_m}{R_a}. \end{aligned}$$

因为所设计的学习系统无需外部教师信号, 是一个高度的自治系统, 所以 OCLA 中的输入符号集合 A 为空集. OCLA 的八元组实际意义见表 2. OCLA 自动机的学习目标是确定一个操作行为序列, 以实现两轮机器人的姿态平衡控制.

表 2 机器人平衡控制问题中 OCLA 八元组实际意义

符号	实际意义
A	外部教师信号
S	机器人倾角和倾角速度, $n_S = 3$
O	电机电压值, $n_O = 3$
Z	机器人控制效果, $n_Z = 2$
R	机器人在状态 s_i 和输入 a_j 条件下依概率 p_{ijk} 实施操作 $o_k (o_k \in O)$
f	机器人状态与控制效果转移方程
ψ	$\psi: S \times A \rightarrow (-1, 0, 1)$
δ	具体描述见式(1)

4 仿真结果与分析

4.1 仿真实验参数设置

OCLA 自动机的仿真参数设置如下: 采样时间为 $t_s = 0.01$ s, 实际状态 θ 的上限值为 $\theta_{\max} = \pi/2$, 下限值为 $\theta_{\min} = -\pi/2$, 输入状态 θ 离散化个数为 $n = 5$ (θ 在 $[-\pi/2 - \pi/3]$ 或 $[\pi/3 \pi/2]$ 时为 s_1 , θ 在 $[-\pi/3 - \pi/6]$ 或 $[\pi/6 \pi/3]$ 时为 s_2 , θ 在 $[-\pi/6 \pi/6]$ 时为 s_3). 实际状态 $\dot{\theta}$ 的上限为 $\dot{\theta}_{\max} = \pi$, 下限为 $\dot{\theta}_{\min} = -\pi$. 将输入状态 θ 离散化的数目等间隔分为 3 段, 初始状态 $\theta = 0.1$ rad, 其他为 0. 用电压表示操作行为集合 $O = \{-5, 0, 5\}$, 每个行为的初始概率为 $p_{ik}(0) = 1/3$, $i = 1, 2, 3$, $k = 1, 2, 3$. 对应初始操作行为为

$$H = - \sum_{k=1}^3 p_{ik} \times \log_2 p_{ik} |_{p_{ik}=\frac{1}{3}} \approx 1.585.$$

此时熵值最大.

当某一行为选择概率值接近于 1 时易发生小概率事件, 从而导致系统的不稳定性增加, 故本文在设计中增加了最优行为选择概率阈值 p_ϵ . 该值大小接近于 1, 当机器人的行为选择概率学习到该值时, 认为已习得某一行为, 从而停止随机选择过程, 以避免小概率事件的发生. 学习率 λ 的选择与学习时间有关, 取值较小时学习时间较长. 当 t 时刻满足 $|\theta| \leq 0.2^\circ$, $|\dot{\theta}| \leq 2^\circ/s$ 且 $p_{ik}(t) > p_\epsilon$ 时, 机器人能够通过学习实现其自主平衡控制, 之后机器人在此状态下持续选择操作 o_k , 直到达到迭代次数 T_f . 仿真中取 $\lambda = 0.05$, $p_\epsilon = 0.94$, $T_f = 1000$.

4.2 结果与分析

经过 20 轮学习后, 系统最终趋于离散状态 s_3 , 图 3 为该状态下映射的 3 个行为选择概率曲线. 初始时刻行为选择概率相同, 机器人状态出现较大振荡. 经过 400 次学习后, 操作行为 o_2 为 0 的概率值逐渐增加, 其他两个行为的概率值相应减小, 机器人此时倾向于选择好的行为 o_2 , 从而使其能够保持在平衡位置. 小概率事件的存在 (在某一时刻机器人选择了小概率的行为) 导致系统状态出现不稳定, 如图 3 中机器人在 420 次训练时选择了小概率行为 o_3 , 从而导致图 4 系统出现暂时的动荡, 但随着学习的进行又趋向于稳定状态. 当检测到机器人 t 时刻满足系统状态 $|\theta| \leq 0.2^\circ$, $|\dot{\theta}| \leq 2^\circ/s$ 和 $p_{32}(t) > p_\epsilon$ 时, 机器人已能通过学习实现自主平衡, 然后按平衡状态下的确定性选择操作 o_2 . 该约束条件的加入避免了小概率事件易导致的破坏性实验.

为了验证本文方法的有效性, 将本文方法与传统的 LQR 方法进行对比仿真研究. 由图 4 可见, 在学习的初始阶段, OCLA 没有任何经验, 振荡较大, 学习速度也较慢, 因此在经过约 430 次学习训练后, OCLA 机器人逐渐学会了平衡. 当机器人达到平衡后, 在第 600 次训练时施加了一个正脉冲干扰信号. 由仿真结果可以看出, OCLA 学习系统能使机器人在约 70 次学习训练后达到平衡, 而 LQR 控制则需要大约 150 次训练次数才能恢复平衡状态. 这说明, 当外界环境发生变化时, OCLA 能使机器人快速地适应环境并作出响应, 具有较好的鲁棒性.

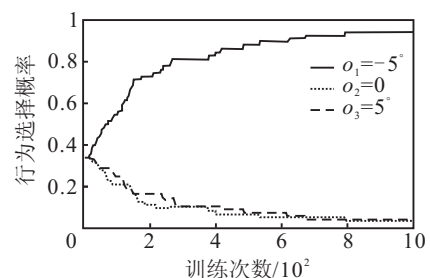


图 3 行为选择概率曲线

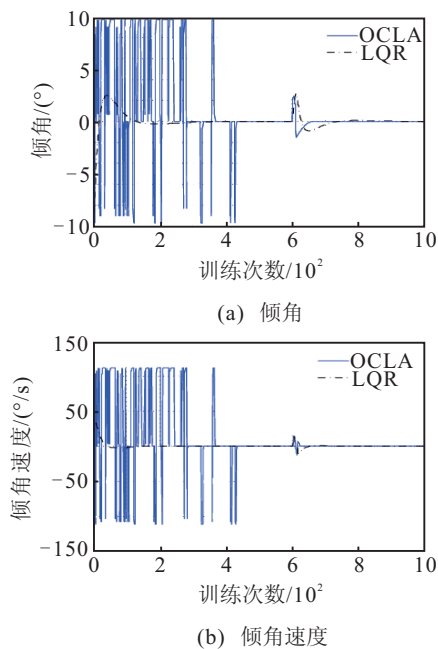


图4 机器人倾角和倾角速度变化曲线

学习结束后各操作行为的概率值为 $p(o_1|s_3) = 0.001$, $p(o_2|s_3) = 0.98$, $p(o_3|s_3) = 0.019$. 代入操作信息熵得到

$$H = - \sum_{k=1}^3 p_{ik} \times \log_2 p_{ik} \approx 0.069, \quad i = 3.$$

图5为机器人操作熵的变化曲线,用熵来反应机器人系统的状态,开始时因其状态出现的概率相等使得此时熵值最大.机器人通过取向信息不断调整行为选择的概率,使机器人慢慢学习不断调整,操作熵的值也逐渐减小并趋于0,最终系统趋于一种稳态,达到控制平衡的目的.

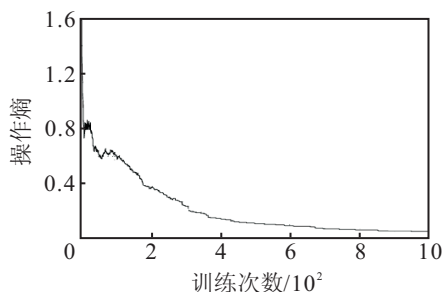


图5 操作熵变化曲线

由仿真结果可以看出,虽然本文所设计的操作条件反射学习自动机模型在初始阶段的控制效果稍差于经典的控制方法,但能有效表达出生物渐近的学习特性,当环境发生变化时具有快速响应性,表现出较好的控制性能.

5 结论

本文在学习自动机理论的框架中基于操作条件反射理论构造了一种仿生学习模型OCLA,并将其应用于两轮机器人的平衡控制.仿真结果表明,OCLA能使机器人自组织地渐近地学习平衡的技能,具有仿生的特性.该方法由于采取状态离散和行为离散的方式,简单易行,但泛化能力较弱,且划分时需要更多的先验知识,对于更为复杂的系统不易得到合理的划分结果,下一步将重点研究如何提高其泛化能力.

参考文献(References)

- [1] Skinner B F. The behavior of organisms[M]. New York: Appleton Century Crofts, 1938: 18-32.
- [2] 任红格, 阮晓钢. 基于Skinner操作条件反射的两轮机器人自平衡控制[J]. 控制理论与应用, 2010, 27(10): 1423-1428.
- [3] 阮晓钢, 陈静. 基于滑模思想和Elman网络的操作条件反射学习控制方法[J]. 控制与决策, 2011, 26(9): 1398-1401.
- [4] Touretzky D S, Saksida L M. Operant conditioning in skinnerbots[J]. Adaptive Behavior, 1997, 5(3/4): 219-247.
- [5] Saksida L M, Raymond S M, Touretzky D S. Shaping robot behavior using principles from instrumental conditioning[J]. Robotics and Autonomous Systems, 1998, 22(3/4): 231-249.
- [6] Gaudio P, Chang C. Adaptive obstacle avoidance with a neural network for operant conditioning: Experiments with real robots[C]. IEEE Int Symposium on Computational Intelligence in Robotics and Automation. New York: IEEE Press, 1997: 13-18.
- [7] Itoh K, Miwa H, Matsumoto M, et al. Behavior model of humanoid robots based on operant conditioning[C]. The 5th IEEE-RAS Int Conf on Humanoid Robots. Tsukuba: Institute of Electrical and Electronic Engineers Computer, 2005: 220-225.
- [8] Pierce D, Kuipers B. Learning to explore and build maps[C]. Proc of the National Conf on Artificial Intelligence. Seattle: AAAI, 1994: 1264-1271.

(下转第939页)