

基于状态聚类的非参数化近似广义策略迭代增强学习算法

季 挺, 张 华[†]

(南昌大学 江西省机器人与焊接自动化重点实验室, 南昌 330031)

摘 要: 为解决当前近似策略迭代增强学习算法普遍存在计算量大、基函数不能完全自动构建的问题, 提出一种基于状态聚类的非参数化近似广义策略迭代增强学习算法(NPAGPI-SC). 该算法利用二级随机采样过程采集样本, 利用 trial-and-error 过程和以样本完全覆盖为目标的估计方法计算逼近器初始参数, 利用 delta 规则和最近邻思想在学习过程中自适应地调整逼近器, 利用贪心策略选择应执行的动作. 一级倒立摆平衡控制的仿真实验结果验证了所提出算法的有效性和鲁棒性.

关键词: 增强学习; 策略迭代; 非参数化; 状态聚类

中图分类号: TP181

文献标志码: A

Nonparametric approximation generalized policy iteration reinforcement learning algorithm based on states clustering

Ji Ting, ZHANG Hua[†]

(Key Lab of Robot & Welding Automation of Jiangxi Province, Nanchang University, Nanchang 330031, China)

Abstract: A nonparametric approximation generalized policy iteration reinforcement learning algorithm based on states clustering(NPAGPI-SC) is proposed to solve the problems such as large calculating quantity and building basis function incompletely automated for the current approximation policy iteration reinforcement learning algorithm. In this algorithm, two stage random sampling process is used to collect samples, the trial-and-error process and the estimation algorithm for covering samples completely are utilized to compute approximator's initial parameters, the delta rule and nearest neighbor method are exploited to adjust the approximator automatically in the learning process, and the greedy strategy is adopted to select an action. The results of simulation on the balancing control of a single inverted pendulum show the effectiveness and robustness of the proposed algorithm.

Keywords: reinforcement learning; policy iteration; nonparametric; states clustering

0 引 言

经典增强学习算法使用表格存储和计算 Q 值(或 V 值), 仅适于解决小规模、离散状态和动作空间中的问题, 实际系统通常工作在大规模、连续状态和动作空间中, 若依然采取表格算法则会导致“维数灾”问题. 近似策略迭代算法是解决“维数灾”问题的主要方法之一, 典型的近似策略迭代增强学习算法 LSPI^[1] 为离线算法, 不仅计算量大, 而且值函数逼近器(包括基函数和参数)依赖先验知识或通过反复试凑确定, 不具备自动构建的能力. online LSPI^[2] 解决了 LSPI 算法的在线计算问题, BLSPI^[3] 在 online LSPI 算法的基础上进一步提高了样本的利用率, 并在一定程度上降低了近似策略迭代算法的计算量, 但这两

种算法依然需要手工构建基函数. 基于核的近似策略迭代增强学习算法^[4-7] 能够在相当程度上解决逼近器的自动构建问题, 但依然需要手动选择基函数宽度、稀疏度阈值等参数.

鉴于此, 本文提出一种基于状态聚类的非参数化近似广义策略迭代增强学习算法(NPAGPI-SC). 该算法能够综合利用离线和在线样本, 不仅计算量较小, 而且只需预先指定增强学习的允许误差率, 便能够在与环境的交互过程中自主构建和调整逼近器的基函数和参数. 将 NPAGPI-SC 应用于一级倒立摆平衡控制问题, 仿真实验结果验证了所提出算法的有效性和鲁棒性, 且相较于 LSPI、online LSPI、BLSPI、KLSPI^[5-6] 等算法, NPAGPI-SC 具有收敛速度更快的

收稿日期: 2016-09-09; 修回日期: 2016-12-13.

基金项目: 国家 863 计划项目(SS2013AA041003).

作者简介: 季挺(1982—), 男, 博士, 从事智能机器人、智能控制的研究; 张华(1964—), 男, 教授, 博士生导师, 从事智能机器人技术、光纤传感、智能金属结构等研究.

[†]通讯作者. E-mail: zhanghua@163.com

优势.

1 NPAGPI-SC算法网络结构

NPAGPI-SC算法网络结构由状态输入层、 Q 值函数逼近器和动作选择器3部分组成,如图1所示.

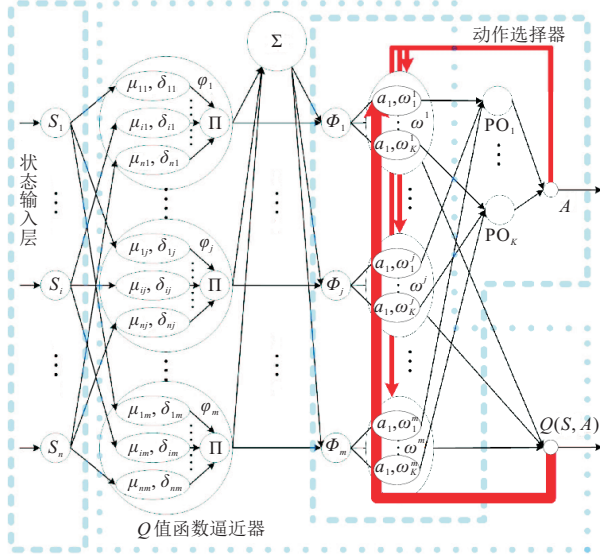


图1 NPAGPI-SC算法网络结构

1.1 状态输入层

状态输入层接收输入状态 $S = (s_1, \dots, s_i, \dots, s_n)^T \in R^n$. 其中: n 为状态空间的维数, s_i 为 S 在各维空间上的分量.

1.2 Q 值函数逼近器

Q 值函数逼近器由基于RBF的线性逼近结构实现,状态动作对 (S, A) 对应的近似 Q 值计算公式为

$$Q(S, A) = \Phi^T(S)\omega(A). \quad (1)$$

$\Phi(S) = (\Phi_1(S), \dots, \Phi_j(S), \dots, \Phi_m(S))^T$ 为状态 S 在各基函数下的归一化隶属度向量,定义为

$$\Phi_j(S) = \frac{\varphi_j(S)}{\sum_{l=1}^m \varphi_l(S)}, \quad j = 1, 2, \dots, m. \quad (2)$$

其中: $\varphi(S) = (\varphi_1(S), \dots, \varphi_j(S), \dots, \varphi_m(S))^T$ 为逼近器的状态基函数向量,其值为状态 S 在各基函数下的隶属度; m 为 $\varphi(S)$ 的维数. $\varphi(S)$ 使用RBF函数定义为

$$\varphi_j(S) = \prod_{i=1}^n \exp \left[-\frac{(s_i - \mu_{ij})^2}{2\delta_{ij}^2} \right], \quad j = 1, 2, \dots, m. \quad (3)$$

其中: $\mu_j = (\mu_{1j}, \dots, \mu_{ij}, \dots, \mu_{nj})^T$ 和 $\delta_j = (\delta_{1j}, \dots, \delta_{ij}, \dots, \delta_{nj})^T$ 分别为基函数 φ_j 的中心和半径, μ_{ij} 和 δ_{ij} 分别为基函数 φ_j 在第 i 维状态分量上的中心和半径.

$\omega = (\omega^1, \dots, \omega^j, \dots, \omega^m)^T$ 为逼近器的插值参

数向量,其意义为在各状态基函数中心执行所选动作 A 的 Q 值,利用插值方法^[8] 定义为

$$\omega^j = \frac{\sum_{k=1}^K d(a_k, A) \cdot \omega_k^j}{\sum_{k=1}^K d(a_k, A)}, \quad j = 1, 2, \dots, m. \quad (4)$$

$$d(a_k, A) = \exp \left(-\frac{(a_k - A)^2}{2\xi^2} \right),$$

$$\xi = \frac{\sum_{k=1}^K |a_k - A|}{K}.$$

其中: K 为可选离散动作的数量; a_k 为第 k 个可选的离散动作; ω_k^j 为逼近器参数,其意义是在状态基函数 φ_j 中心 μ_j 执行动作 a_k 的投票数,采用 delta 规则进行更新,定义为

$$\omega_k^j(t+1) = \omega_k^j(t) + \eta \Delta \Phi_j(S_t) \frac{d(a_k, A)}{\sum_{k=1}^K d(a_k, A)}. \quad (5)$$

这里: η 为 Q 值函数逼近器平均学习率,由初始状态基函数构建过程确定; Δ 为 TD 误差,定义为

$$\Delta = r + \gamma \max_{A_{t+1}} Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t), \quad (6)$$

r 为即时奖励, γ 为折扣率.

状态基函数向量 φ (包括其参数和维数) 与插值参数向量 ω 的维数由 Q 值函数逼近器自动构建过程 (见第2节) 依据采样样本确定,并在学习过程中自适应调整.

1.3 动作选择器

动作选择器采用贪心策略实现,计算公式为

$$A = \arg \max_{a_k} (PO_1, \dots, PO_k, \dots, PO_K), \quad (7)$$

其中 PO_k 是对当前状态 S 下执行动作 a_k 的投票结果,定义为

$$PO_k = \sum_{j=1}^m \Phi_j(S) \omega_k^j, \quad (8)$$

$$E(PO_k) = \sum_{j=1}^m \Phi_j(S) E(\omega_k^j), \quad (9)$$

$$\text{Var}(PO_k) = \sum_{j=1}^m \Phi_j^2(S) \text{Var}(\omega_k^j). \quad (10)$$

由式(9)和(10)可以看出,投票器算法可有效降低单个 ω_k^j 估计误差对动作选择结果的影响,能够大大提升算法对动作噪声的抗干扰能力.

2 Q 值函数逼近器自动构建过程

Q 值函数逼近器自动构建过程由二级随机采样过程、核心状态聚类生成过程、初始状态基函数构建

过程组成.

2.1 二级随机采样过程

近似增强学习方法中的近似不仅指表示近似,也指样本近似. 在大规模、连续状态空间中遍历每个状态是不可能的,因此需要一种好的采样方法尽量保证样本空间对状态空间的覆盖率. 本文提出一种二级随机采样过程,该过程分为局部随机采样过程和全局随机采样过程两个层次,依据预先指定的增强学习允许误差率 ε 自动收集状态样本.

2.1.1 局部随机采样过程

局部随机采样过程如下.

Step 1: 采用随机策略运行一次目标系统. 设置采样次数 $tm = 1$, 目标系统运行次数 $tr = 1$, 采样缓存长度 $L = step$, $step$ 为目标系统本次运行持续的时间步数(同时也是目标系统本次运行采集样本的数量); 将本次采样结果记为 sam_{tm} , 通过下式分别计算均值 $mean_{tm}$ 及其在各维分量上的平均绝对离差向量 mad_{tm} :

$$mean_{tm} = \frac{\sum_{h=1}^{ns} (sam_{tm})_h}{ns}, \quad (11)$$

$$mad_{tm} = \frac{\sum_{h=1}^{ns} \text{abs}((sam_{tm})_h - mean_{tm})}{ns}, \quad (12)$$

其中 ns 为 sam_{tm} 中样本的数量.

Step 2: 开始单次采样. $tm = tm + 1$, 初始化单次采样新增非重复样本数量 $ln = 0$, 使 $temp = sam_{tm-1}$.

Step 3: 采用随机策略运行一次目标系统. $tr = tr + 1$, 收集本次运行采集到的样本集合 T_{sam} , 同时更新 $step$ 和 rep 的值, rep 为 T_{sam} 与 $temp$ 中重复样本的数量.

Step 4: 使得 $temp = temp \cup T_{sam}$, 更新 ln 与 L 的值, 有

$$ln = ln + step - rep. \quad (13)$$

$$L = \begin{cases} L - 1, & rep == step; \\ \text{round}\left(\frac{L(tr - 1) + step}{tr}\right), & \text{otherwise.} \end{cases} \quad (14)$$

若 $ln < L$, 则执行Step 3, 否则, 令 $sam_{tm} = temp$, 执行Step 5.

Step 5: 由式(11)和(12)计算 $mean_{tm}$ 和 mad_{tm} .

Step 6: 计算 sam_{tm} 相对 sam_{tm-1} 的变化率 $rate$, 有

$$rate = \|(rate_1, \dots, rate_i, \dots, rate_n)\|_2. \quad (15)$$

$$rate_i = \begin{cases} 1, & bc_i \leq fl_i | bl_i \leq fc_i; \\ \frac{|fc_i - fl_i| + |bc_i - bl_i|}{\max(fbcl) - \min(fbcl)}, & \text{otherwise.} \end{cases}$$

$$fbcl = \{fc_i, bc_i, fl_i, bl_i\}.$$

$$fc = mean_{tm} - mad_{tm}.$$

$$bc = mean_{tm} + mad_{tm}.$$

$$fl = mean_{tm-1} - mad_{tm-1}.$$

$$bl = mean_{tm-1} + mad_{tm-1}.$$

Step 7: 若 $rate < \varepsilon$, 则局部采样过程结束, 否则执行Step 2.

2.1.2 全局随机采样过程

全局随机采样过程如下.

Step 1: 目标系统运行一次局部随机采样过程. 设置局部采样次数 $TM = 1$, 将本次局部采样结果记为 SAM_{TM} , 类似式(11)和(12)分别计算均值 $MEAN_{TM}$ 及其在各维分量上的平均绝对离差向量 MAD_{TM} .

Step 2: 目标系统运行一次局部随机采样过程, 得到本次局部采样结果 T_{SAM} . $TM = TM + 1$, $SAM_{TM} = SAM_{TM-1} \cup T_{SAM}$.

Step 3: 类似式(11)和(12)计算 $MEAN_{TM}$ 和 MAD_{TM} .

Step 4: 类似式(15)计算 SAM_{TM} 相对 SAM_{TM-1} 的变化率 $RATE$.

Step 5: 若 $RATE < \varepsilon$, 最终样本集 $SAM = SAM_{TM}$, 则全局采样过程结束, 否则执行Step 2.

2.2 核心状态聚类生成过程

核心状态聚类生成过程的目的是为初始状态基函数构建过程提供初值, 包括trial-and-error过程和总体生成过程两个层次.

2.2.1 trial-and-error过程

trial-and-error过程^[9]是一种使用相对准则判定样本集最佳聚类数的方法, 本文基于轮廓指标^[10]、采用 K 均值聚类算法实现该过程, 具体方法如下.

Step 1: 初始化输入样本集 sam 的可能聚类结果集合 $clus = \emptyset$ 、 $clus$ 对应的有效性集合 $val = \emptyset$; 使用式(16)初始化 sam 的可能聚类数集合 num , num 的最大取值为 \sqrt{N} ^[11], N 为 sam 中样本的数量. 设置每种聚类数下运行 K 均值聚类算法的次数 RTM (一般3~10次), 有

$$num = \begin{cases} \{1, 2, \dots, [\sqrt{N}]\}, & num \text{ is } \emptyset; \\ num, & \text{otherwise.} \end{cases} \quad (16)$$

Step 2: 遍历 num , 对当前遍历的聚类数 num_v 运

行 K 均值聚类算法 RTM 次, 得到聚类结果 $(TC_1, \dots, TC_w, \dots, TC_{RTM})$ 及其相应的聚类有效性 $(TV_1, \dots, TV_w, \dots, TV_{RTM})$, TV_w 采用轮廓指标^[10]计算, 有

$$TV_w = \frac{1}{\text{num}_v} \sum_{p=1}^{\text{num}_v} \frac{1}{n_p} \sum_{q: x_q \in (TC_w)_p} \frac{b_q - a_q}{\max(b_q, a_q)}. \quad (17)$$

$$a_q = d_{\text{avg}}^{\text{ps}}(x_q, (TC_w)_p - \{x_q\}),$$

$$b_q = \min_{u=1,2,\dots,\text{num}_v, u \neq p} d_{\text{avg}}^{\text{ps}}(x_q, (TC_w)_u).$$

其中: n_p 为聚类 $(TC_w)_p$ 的基数, a_q 为 $(TC_w)_p$ 中元素 x_q 到 $(TC_w)_p$ 中其他元素的平均距离, b_q 为 x_q 到除 $(TC_w)_p$ 之外最接近 x_q 的聚类的平均距离, $d_{\text{avg}}^{\text{ps}}(\cdot, \cdot)$ 为一个点到一个集合的平均距离测度.

Step 3: 将 num_v 下的最优聚类及其有效性存入 clus_v 和 val_v , 有

$$id = \arg \max_w (TV_1, \dots, TV_w, \dots, TV_{RTM}),$$

$$\text{clus}_v = TC_{id}; \quad (18)$$

$$\text{val}_v = \frac{\sum_{w=1}^{\text{RTM}} TV_w}{\text{RTM}}. \quad (19)$$

Step 4: 若 num 遍历结束, 则输出 sam 的最佳聚类数 Num 和聚类结果 clu , trial-and-error 过程结束, 否则执行 Step 2, 有

$$ID = \arg \max_v (\text{val}),$$

$$\text{Num} = \text{num}_{ID}; \quad (20)$$

$$\text{clu} = \text{clus}_{ID}. \quad (21)$$

2.2.2 总体生成过程

核心状态聚类总体生成过程如下.

Step 1: 初始化全局样本集 SAM 的可能聚类数集合 $\text{cnum} = \emptyset$.

Step 2: 目标系统运行一次局部随机采样过程, 得到样本集 T_{SAM} .

Step 3: 在 T_{SAM} 上运行 trial-and-error 过程, 获取其最佳聚类数 NUM , 并更新 cnum , 有

$$\text{cnum} = \text{cnum} \cup \text{NUM}. \quad (22)$$

Step 4: 若全局采样过程结束, 则初始化 SAM 的可能聚类数集合为 cnum , 在 SAM 上运行 trial-and-error 过程, 获取其最佳聚类结果 CLU, 并将 CLU 作为核心状态聚类, 否则执行 Step 2.

Step 5: 计算 CLU 中每个聚类的中心和半径, 并将

其存入核心聚类中心矩阵 cen 和半径矩阵 rad . CLU 中单个聚类 C 的中心 C_c 为 C 中所有元素的均值, 半径 C_r 为 C 中每个元素与其中心 C_c 在各维分量上的距离均值向量, 分别为

$$C_c = \frac{\sum_{z=1}^{\text{nc}} C_z}{\text{nc}}, \quad (23)$$

$$C_r = \frac{\sum_{z=1}^{\text{nc}} \text{abs}(C_z - C_c)}{\text{nc}}. \quad (24)$$

其中: C_z 为聚类 C 中的一个元素, nc 为聚类 C 中元素的数量.

2.3 初始状态基函数构建过程

状态基函数如式 (3) 所示. 为了提升状态基函数对样本空间的覆盖率、使其尽可能完整地反映状态空间的分布, 本文提出一种以样本空间完全覆盖为目标估计方法, 通过采集到的样本计算逼近器的状态基函数和平均学习率, 具体过程如下.

Step 1: 初始化基函数向量 φ .

Step 2: 遍历样本集 SAM. 通过下式判定与当前遍历样本 C_{sam} 最邻近的基函数 φ_{near} :

$$\varphi_{\text{near}} = \arg \min_{\varphi_j} \{|C_{\text{sam}} - \mu_j| | j = 1, 2, \dots, m\}; \quad (25)$$

通过下式判定 C_{sam} 是否落在 φ_{near} 内部:

$$\text{sum}(\text{abs}(C_{\text{sam}} - \mu_{\text{near}}) \leq \delta_{\text{near}}) == n. \quad (26)$$

若是, 则继续遍历过程 Step 2, 否则结束本次遍历, 执行 Step 3. 若 SAM 中所有样本均落在最邻近基函数内部, 则表明状态基函数完成了对样本空间的完全覆盖, 执行 Step 5. 其中, μ_{near} 和 δ_{near} 分别为 φ_{near} 的中心和半径.

Step 3: 新增一个状态基函数, 将其中心初始化为 C_{sam} , 并加入 φ 中.

Step 4: 通过式 (27) 和 (28) 调整 φ 中所有基函数的中心和半径, 执行 Step 2, 有

$$\mu_j = \mu_j + \frac{\sum_{h=1}^{\text{NS}} \text{be}_{hj} \cdot (\text{SAM}_h - \mu_j)}{\sum_{h=1}^{\text{NS}} \text{be}_{hj}}. \quad (27)$$

$$\delta_j = \frac{\sum_{h=1}^{\text{NS}} \text{be}_{hj} \cdot |\text{SAM}_h - \mu_j|}{\sum_{h=1}^{\text{NS}} \text{be}_{hj}}. \quad (28)$$

$$be_{hj} = \begin{cases} 1, & \text{SAM}_h == \mu_j; \\ \frac{1}{\sum_{l=1}^m \frac{1}{|SAM_h - \mu_l|}}, & \text{otherwise.} \end{cases}$$

其中: NS 为 SAM 中样本的数量, be_{hj} 为样本 SAM_h 到状态基 φ_j 的归一化隶属度.

Step 5: 通过式 (29) 计算逼近器的平均学习率 η , 初始状态基函数构建过程结束, 有

$$\eta = \frac{2}{\bar{X} + M_e(X)}. \quad (29)$$

其中 $X = (X_1, \dots, X_j, \dots, X_m)^T$ 为各基函数的构建样本数向量, 表示为

$$X_j = \sum_{h=1}^{NS} be_{hj}, \quad j = 1, 2, \dots, m.$$

3 Q 值函数逼近器自适应调整过程

Q 值函数逼近器自适应调整过程包括逼近器状态基函数和逼近器参数的自适应调整. 其中, 参数的自适应调整如式 (5) 所示, 状态基函数的自适应调整过程如下所示.

Step 1: 若输入状态 S 属于状态基函数的构建样本集 SAM, 则继续执行 Step 1, 判定下一输入状态, 否则将 S 加入构建样本集 SAM, 执行 Step 2.

Step 2: 通过式 (30) 搜索与 S 最邻近的基函数 φ_{near} , 有

$$\varphi_{near} = \arg \lim_{\varphi_j} \{S - \mu_j | j = 1, 2, \dots, m\}. \quad (30)$$

Step 3: 若 φ_{near} 为初始状态基函数, TD 误差率 $R_\Delta > \varepsilon$, 通过下式判定 S 落在 φ_{near} 外部:

$$R_\Delta = \frac{|\Delta|}{\max(\text{abs}(r))}, \quad (31)$$

$$\text{sum}(\text{abs}(S - \mu_j) \leq \delta_j) == n, \quad (32)$$

则新增状态基函数 φ_{new} , 并初始化其中心 $\mu_{new} = S$, 半径 $\delta_{new} = \delta_{near}$, 执行 Step 4, 否则直接执行 Step 4.

Step 4: 采用 delta 规则调整状态基函数的中心和半径, 对于 φ_j 有

$$\mu_j = \mu_j + \eta \cdot R_\Delta \cdot BE_j \cdot (S - \mu_j), \quad (33)$$

$$\delta_j = \delta_j + \eta \cdot R_\Delta \cdot BE_j \cdot (|S - \mu_j| - \delta_j). \quad (34)$$

其中

$$BE_j = \frac{\exp\left[-\frac{(S - \mu_j)^2}{2\delta_j^2}\right]}{\sum_{l=1}^m \exp\left[-\frac{(S - \mu_l)^2}{2\delta_l^2}\right]}.$$

本次调整结束后执行 Step 1, 开始下一次状态基函数调整过程.

4 NPAGPI-SC 算法步骤

NPAGPI-SC 算法步骤分为算法网络结构构建和自主学习两个阶段.

4.1 算法网络结构构建阶段

算法网络结构构建阶段如图 2 所示, 具体步骤如下所示.

Step 1: 指定增强学习的允许误差率 ε .

Step 2: 使用全局和局部二级随机采样过程采集误差低于 ε 的状态样本集合 SAM. 在每次局部采样结束后利用 trail-and-error 过程计算局部样本集合 TSAM 的最佳聚类数 NUM, 并使用式 (22) 得到 SAM 的可能最佳聚类数集合 cnum.

Step 3: 依据 cnum, 使用 trail-and-error 过程计算 SAM 的核心状态聚类 CLU. 使用式 (23) 和 (24) 分别计算 CLU 的中心矩阵 cen 和半径矩阵 rad.

Step 4: 分别利用 cen 和 rad 初始化基函数向量 φ 的中心和半径, 使用初始状态基函数构建过程计算 φ 和逼近器平均学习率 η , 将 Q 值函数逼近器参数向量 $(\omega_1^1, \dots, \omega_K^1, \dots, \omega_1^j, \dots, \omega_K^j, \dots, \omega_1^m, \dots, \omega_K^m)^T$ 初始化为 0 向量.

Step 5: 依据 φ 、 η 、逼近器初始参数向量构建 NPAGPI-SC 算法的网络结构 (如图 2 所示).

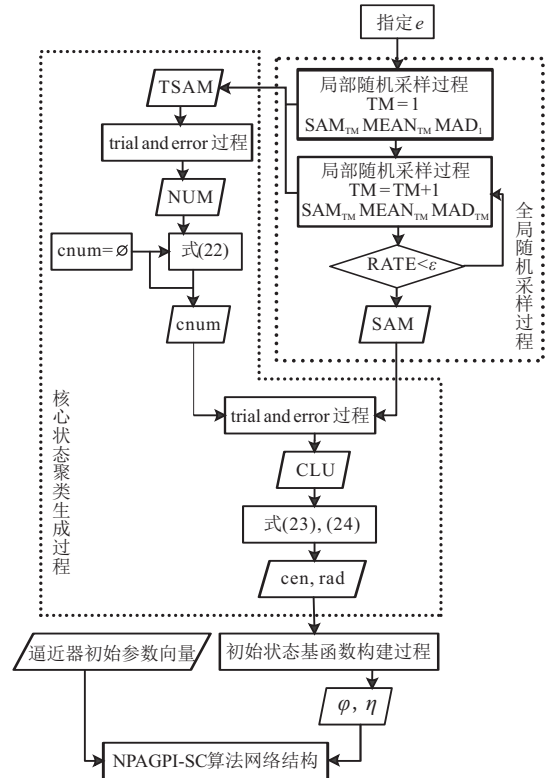


图 2 NPAGPI-SC 算法网络结构构建阶段

4.2 自主学习阶段

NPAGPI-SC 算法网络结构构建完毕后将进入自主学习阶段, 该阶段采用广义策略迭代思想实现, 策

略评估在未收敛时便与策略改进交替进行,具体步骤如下所示.

Step 1: 目标系统开始一次尝试. 初始化当前状态 $S_t = S_0$. 若本次尝试未结束,则循环执行下列步骤:

Step 1.1: 使用式(3)计算 S_t 在各状态基函数的隶属度向量 $\varphi(S_t)$;

Step 1.2: 使用式(2)计算 S_t 的归一化状态基函数隶属度向量 $\Phi(S_t)$;

Step 1.3: 使用式(8)进行策略评估,使用式(7)选择当前状态下应采取的动作 A_t , 获取即时奖励 r , 且环境转移到下一状态 S_{t+1} ;

Step 1.4: 使用式(6)计算TD误差 Δ ;

Step 1.5: 依据式(5)调整 Q 值函数逼近器参数向量进行策略改进,依据 Q 值函数逼近器自适应调整过程更新构建样本集和状态基函数. 设置当前状态 $S_t = S_{t+1}$.

Step 2: 若本次尝试成功则学习算法结束,否则执行Step 1再次进行尝试.

从以上步骤可以看出, NPAGPI-SC 在学习阶段的计算复杂度主要由策略评估复杂度 $O(Km)$ 、策略选择复杂度 $O(K)$ 、策略改进复杂度 $O(K + Km + K^2m)$ 和基函数调整复杂度 $O(m)$ 组成, 即 $O(2K + (K^2 + 2K + 1)m)$, 远小于其他近似策略迭代增强学习算法 $O(m^3)$ 的时间复杂度, 可以更好地满足在线计算的要求.

5 一级倒立摆平衡控制仿真实验

5.1 实验描述

一级倒立摆平衡控制是增强学习研究的标准测试问题, 本文通过对其进行仿真实验来验证 NPAGPI-SC 算法的有效性和鲁棒性. 为了方便与已有算法进行对比, 采用文献[1]的方法建立目标系统模型

$$\ddot{\theta} = \frac{g \cdot \sin \theta - \frac{\cos \theta \cdot (F + MP \cdot \text{len} \cdot \dot{\theta}^2 \cdot \sin \theta)}{MC + MP}}{\text{len} \cdot \left(\frac{4}{3} - \frac{MC \cdot \cos^2 \theta}{MC + MP} \right)} \quad (35)$$

其中: 系统的状态 $S = [\theta \ \dot{\theta}]^T$, θ 为摆杆偏离垂直方向的角度, $\dot{\theta}$ 为角速度; F 为作用于小车质心上的力, 最大值为 50 N, 向右取 +, 向左取 -; MC 为小车质量, $MC = 8.0 \text{ kg}$; MP 为摆杆质量, $MP = 2.0 \text{ kg}$; len 为摆杆长度的一半, $\text{len} = 0.5 \text{ m}$; g 为重力加速度, $g = 9.8 \text{ m/s}^2$.

实验中, 仿真时间步设置为 0.1 s, 折扣率 $\gamma = 0.95$. 若倒立摆在一次尝试中保持指定时间步不倒, 则认为本次实验成功; 否则, 若摆杆偏离垂直方向角

度超过 90° , 则认为本次实验失败, 相应的即时奖励 r 如下所示:

$$r = \begin{cases} 0, & |\theta| \leq 90^\circ; \\ -1, & \text{otherwise.} \end{cases} \quad (36)$$

5.2 实验结果

5.2.1 NPAGPI-SC 取不同参数的实验结果

设置实验目标为倒立摆运行 3 000 时间步不倒, 实验最大尝试次数为 300, 当 $\varepsilon = 0.1$ 时, 离散动作集合分别取 2 个离散值 $\{-50 \text{ N}, +50 \text{ N}\}$ 、3 个离散值 $\{-50 \text{ N}, 0 \text{ N}, +50 \text{ N}\}$ 、4 个离散值 $\{-50 \text{ N}, -25 \text{ N}, +25 \text{ N}, +50 \text{ N}\}$ 、5 个离散值 $\{-50 \text{ N}, -25 \text{ N}, 0 \text{ N}, +25 \text{ N}, +50 \text{ N}\}$ 、6 个离散值 $\{-50 \text{ N}, -33 \text{ N}, -16 \text{ N}, +16 \text{ N}, +33 \text{ N}, +50 \text{ N}\}$ 、7 个离散值 $\{-50 \text{ N}, -33 \text{ N}, -16 \text{ N}, 0 \text{ N}, +16 \text{ N}, +33 \text{ N}, +50 \text{ N}\}$ 、8 个离散值 $\{-50 \text{ N}, -37.5 \text{ N}, -25 \text{ N}, -12.5 \text{ N}, +12.5 \text{ N}, +25 \text{ N}, +37.5 \text{ N}, +50 \text{ N}\}$ 、9 个离散值 $\{-50 \text{ N}, -37.5 \text{ N}, -25 \text{ N}, -12.5 \text{ N}, 0 \text{ N}, +12.5 \text{ N}, +25 \text{ N}, +37.5 \text{ N}, +50 \text{ N}\}$ 、10 个离散值 $\{-50 \text{ N}, -40 \text{ N}, -30 \text{ N}, -20 \text{ N}, -10 \text{ N}, +10 \text{ N}, +20 \text{ N}, +30 \text{ N}, +40 \text{ N}, +50 \text{ N}\}$ 、11 个离散值 $\{-50 \text{ N}, -40 \text{ N}, -30 \text{ N}, -20 \text{ N}, -10 \text{ N}, 0 \text{ N}, +10 \text{ N}, +20 \text{ N}, +30 \text{ N}, +40 \text{ N}, +50 \text{ N}\}$, 在以上离散动作集合取值下分别进行 100 次独立仿真运算, 实验结果如表 1 所示.

表 1 离散动作集合在不同取值下的实验结果

离散动作集合	成功率/%	最小尝试次数	最大尝试次数	平均尝试次数	平均样本数量	平均基函数数量(学习前/后)
2个离散值	100	3	76	7.65	167.69	11.92/12.07
3个离散值	100	3	81	11.92	214.61	10.85/11.01
4个离散值	100	3	77	11.61	258.96	11.62/11.75
5个离散值	100	3	85	10.14	256.24	10.74/10.99
6个离散值	100	3	141	9.97	279.65	11.36/11.55
7个离散值	100	3	117	8.25	272.74	10.96/11.19
8个离散值	100	3	131	9.62	271.37	11.26/11.42
9个离散值	100	3	218	11.97	306.39	11.28/11.55
10个离散值	100	3	156	10.95	284.52	11.67/11.91
11个离散值	100	3	60	5.96	285.51	11.16/11.32

由表 1 可见, 不论离散动作集合如何取值, NPAGPI-SC 算法均能在较短尝试次数内学习到一级倒立摆系统的平衡控制策略, 验证了算法在离散动作集合不同取值下的有效性. 若每个被选择执行的离散动作都叠加 $[-10 \text{ N}, +10 \text{ N}]$ 的均匀噪声, 则在离散动作集合的每种取值下分别进行 100 次独立仿真运算, 实验结果如表 2 所示.

对比表 1 和表 2 可以看出, 当离散动作叠加噪声后, 离散动作集合在不同取值下的采样样本数量和基函数数量均略有增加, 但并未影响算法的成功率, 且平均尝试次数并未明显增加, 显示出算法在离散动作集合不同取值下均具有较好的鲁棒性.

表2 离散动作集合在不同取值下叠加噪声后的实验结果

离散动作集合	成功率/%	最小尝试次数	最大尝试次数	平均尝试次数	平均样本数量	平均基函数数量(学习前/后)
2个离散值	100	3	63	17.82	300.92	13.1/13.5
3个离散值	100	3	69	11.42	290.21	13.29/13.69
4个离散值	100	3	72	13.4	310.07	13.02/13.37
5个离散值	100	3	83	9.43	307.1	12.79/13.08
6个离散值	100	3	90	12.31	313.09	12.47/12.94
7个离散值	100	3	133	10.4	336.27	12.96/13.33
8个离散值	100	3	92	10.06	294.18	12.39/12.81
9个离散值	100	3	47	8.18	330.02	12.6/12.84
10个离散值	100	3	77	8.98	307.39	12.05/12.34
11个离散值	100	3	74	10.05	326.6	12.39/12.71

设置实验目标为倒立摆运行3000时间步不倒,实验最大尝试次数为300,当离散动作集合取3个离散值 $\{-50N, 0N, +50N\}$, ϵ 取0.5、0.1、0.05、0.01分别进行100次独立仿真运算,实验结果如表3所示。

表3 ϵ 在不同取值下的实验结果

ϵ 取值	成功率/%	最小尝试次数	最大尝试次数	平均尝试次数	平均样本数量	平均基函数灵敏度(学习前/后)
0.5	100	3	42	6.18	68.88	8.88/9.09
0.1	100	3	81	11.92	214.61	10.85/11.01
0.05	100	3	101	8.62	396.99	12.2/12.33
0.01	100	3	193	14.19	1715.49	16.95/16.98

由表3可见,随着 ϵ 精度的提高,采样样本数量和基函数数量均有所增加,但不论 ϵ 如何取值,NPAGPI-SC算法均能在较短尝试次数内学习到一级倒立摆系统的平衡控制策略,验证了算法在 ϵ 不同取值下的有效性。若每个被选择执行的离散动作都叠加 $[-10N, +10N]$ 的均匀噪声,则在不同 ϵ 取值下分别进行100次独立仿真运算,实验结果如表4所示。

表4 动作叠加噪声后 ϵ 在不同取值下的实验结果

ϵ 取值	成功率/%	最小尝试次数	最大尝试次数	平均尝试次数	平均样本数量	平均基函数灵敏度(学习前/后)
0.5	100	3	52	9.51	90.46	10.28/11.21
0.1	100	3	69	11.42	290.21	13.29/13.69
0.05	100	3	96	16.31	591.76	15.02/15.37
0.01	100	3	192	17.76	2880.54	19.51/19.54

对比表3和表4可以看出,当在离散动作上叠加噪声后, ϵ 在不同取值下的采样样本数量和基函数数量均有所增加,但并未影响算法的成功率,且平均尝试次数并未有明显增加,显示出算法在 ϵ 不同取值下均具有较好的鲁棒性。

设置实验目标为倒立摆运行3000时间步不倒,当离散动作集合取3个离散值 $\{-50N, 0N, +50N\}$, ϵ 取0.1时进行100次独立仿真运算,学习成功的尝试次数分布如图3所示,典型摆杆角度和角速度的学习过程曲线如图4所示。

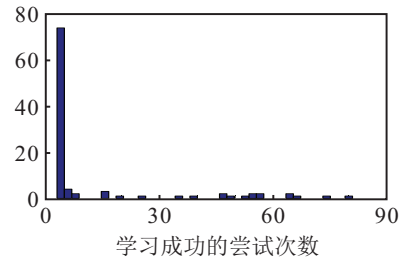
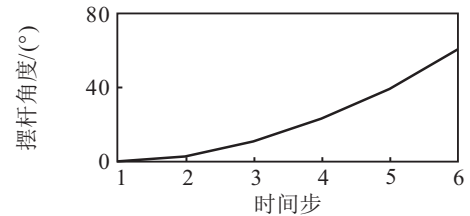
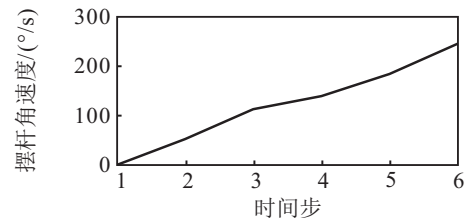


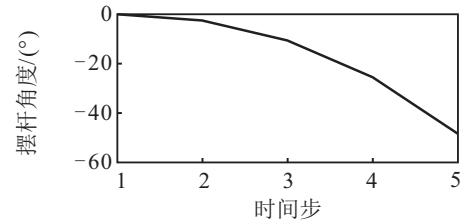
图3 进行100次仿真运算尝试次数的分布



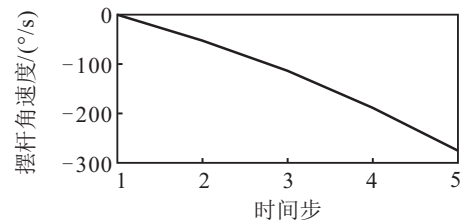
(a) 摆杆角度(第1次尝试)



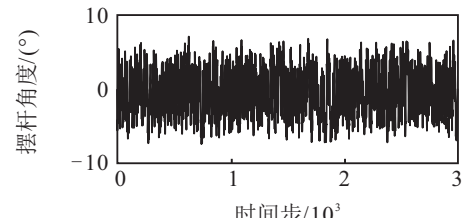
(b) 摆杆角速度(第1次尝试)



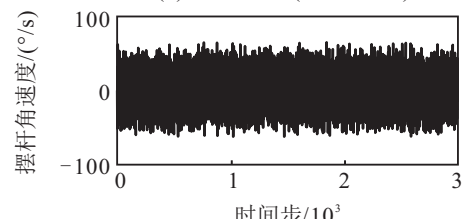
(c) 摆杆角度(第2次尝试)



(d) 摆杆角速度(第2次尝试)



(e) 摆杆角度(第3次尝试)



(f) 摆杆角速度(第3次尝试)

图4 典型学习过程曲线

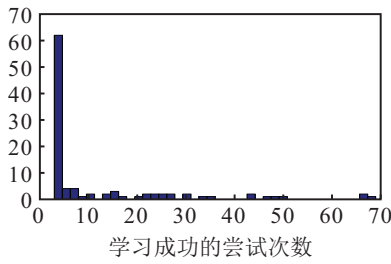


图5 100次独立仿真运算尝试次数的分布(有噪声下)

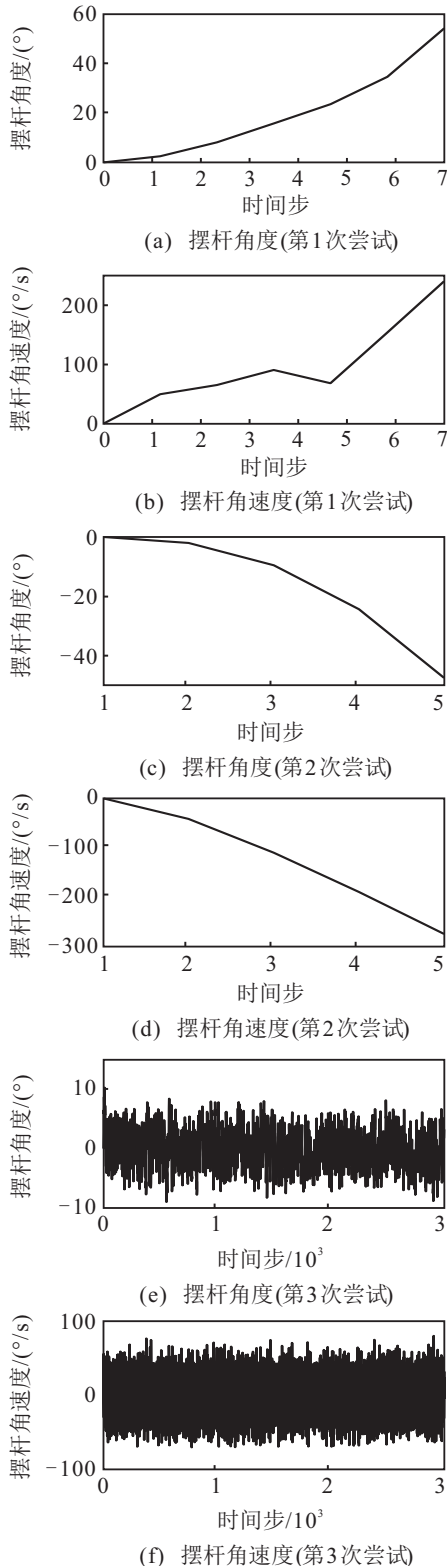


图6 典型学习过程曲线(有噪声下)

由图3可见,算法基本可在第3次尝试时得到正确的控制策略.由图4可见,学习成功后摆杆的摆动角度基本稳定在 $[-8^\circ, +8^\circ]$,角速度基本稳定在 $[-60^\circ/\text{s}, +60^\circ/\text{s}]$.以上数据表明算法的学习速度和控制精度都较为理想.

若每个被选择执行的离散动作都叠加 $[-10\text{N}, +10\text{N}]$ 的均匀噪声,则进行100次独立仿真运算,学习成功的尝试次数分布如图5所示,典型摆杆角度和角速度的学习过程曲线如图6所示.

由图5可见,算法基本可在第3次尝试时得到正确的控制策略.由图6可见,学习成功后摆杆的摆动角度基本稳定在 $[-10^\circ, +10^\circ]$,摆杆角速度基本稳定在 $[-80^\circ/\text{s}, +80^\circ/\text{s}]$.通过对离散动作叠加噪声前后的典型学习过程曲线进行比较可以看出,噪声对NPAGPI-SC算法学习速度和控制精度未造成较大影响,显示出算法在控制性能方面具有较强的鲁棒性.

5.2.2 与典型近似策略迭代增强学习算法的对比

本节依旧采用一级倒立摆平衡控制仿真实验比较各类近似策略迭代增强学习算法的性能.为了与文献[3]和文献[6]的实验条件保持一致,使离散动作集合取3个离散值 $\{-50\text{N}, 0\text{N}, +50\text{N}\}$,且每个被选择执行的离散动作都叠加 $[-10\text{N}, +10\text{N}]$ 的均匀噪声.

online LSPI、BLSPI等在线近似策略迭代增强学习算法一般通过不同参数条件下获得控制策略的速度进行评价.设置实验目标为倒立摆运行3000时间步不倒,将 ϵ 分别为0.5、0.1、0.05、0.01时, NPAGPI-SC进行100次独立仿真运算的结果与onlineLSPI、BLSPI算法的最优实验结果^[3]进行比较,3种算法在不同尝试次数下的平均平衡时间步数如表5所示.

表5 与典型在线近似策略迭代增强学习算法的对比

算 法	50次 尝试	100次 尝试	150次 尝试	200次 尝试	250次 尝试	300次 尝试
NPAGPI-SC $\epsilon = 0.5$	2980	3000	3000	3000	3000	3000
NPAGPI-SC $\epsilon = 0.1$	2944	3000	3000	3000	3000	3000
NPAGPI-SC $\epsilon = 0.05$	2841	3000	3000	3000	3000	3000
NPAGPI-SC $\epsilon = 0.01$	2775	2820	2955	3000	3000	3000
online LSPI	27	35	36	37	928	2112
BLSPI LSTD-Q	921	1800	2187	2576	2771	3000
BLSPI LSPE-Q	34	1532	2701	3000	3000	3000

由表5可见,随着 ϵ 精度的提高, NPAGPI-SC算法的学习速度逐渐降低,但总体性能要优于online LSPI和BLSPI.

LSPI与KLSPI等离线近似策略迭代增强学习算法一般通过在不同样本数目条件下策略迭代获得控制策略的性能进行评价^[6].从实现方法上看,LSPI和KLSPI在每次迭代时都将样本数据累加到矩阵中,其

一次迭代过程所需的样本数量大致与NPAGPI-SC一次尝试相当,因此NPAGPI-SC的平均尝试次数基本等价于LSPI与KLSPI的平均迭代次数.当 $\varepsilon = 0.1$ 时,NPAGPI-SC独立运行100次的平均采样数量为290个,与文献[6]中50周期(约300个样本)的样本数量相当.选择与文献[6]相同的实验方法,即策略性能评价的最大运行实际步长为2000、每个策略独立运行1000次,得到的实验结果对比如表6所示.

表6 采用50周期训练样本获得的策略性能

算 法	平均迭代 次数	最小 时间步	最大 时间步	平均 时间步
NPAGPI-SC $\varepsilon = 0.1$	9.32	2000	2000	2000
LSPI	20	6	2000	460
KLSPI	20	898	2000	1997.9

表6中NPAGPI-SC的平均时间步由100次独立仿真结果中随机抽取10个所学策略分别运行1000次得到的均值计算.当 $\varepsilon = 0.05$ 时,NPAGPI-SC独立运行100次的平均采样数量为592个,与文献[6]中100周期(约600个样本)的样本数量相当.选择与文献[6]相同的实验方法,得到的实验结果对比如表7所示.

表7 采用100周期的训练样本获得的策略性能

算 法	平均迭代 次数	最小 时间步	最大 时间步	平均 时间步
NPAGPI-SC $\varepsilon = 0.05$	11.95	2000	2000	2000
LSPI	20	6	2000	896
KLSPI	20	2000	2000	2000

实验数据表明,无论采用50周期的训练样本还是100周期的训练样本,NPAGPI-SC算法都可以在更少的平均迭代次数下达到或超过LSPI和KLSPI算法的策略性能.

6 结 论

本文实现了一种基于状态聚类的非参数化近似广义策略迭代增强学习算法NPAGPI-SC.相较典型的近似策略迭代增强学习算法,NPAGPI-SC不仅计算复杂度小,且只需指定增强学习的允许误差率便能完全自主地构建、调整基函数和参数,学习到目标系统的控制策略.一级倒立摆平衡控制的仿真实验结果验证了所提出的算法在不同参数下的有效性和鲁棒性,相较于目前典型的近似策略迭代增强学习算

法,NPAGPI-SC具有收敛速度更快的优势.

参考文献(References)

- [1] Lagoudakis Michail G, Parr Ronald. Least squares policy iteration[J]. J of Machine Learning Research, 2003, 4(6): 1107-1149.
- [2] Busoniu Lucian, Ernst Damien, Schutter Bart De, et al. Online least-squares policy iteration for reinforcement learning control[C]. Proc of the 2010 American Control Conf. Baltimore: IEEE Press, 2010: 486-491.
- [3] 周鑫, 刘全, 傅启明, 等. 一种批量最小二乘策略迭代方法[J]. 计算机科学, 2014, 41(9): 232-238. (Zhou X, Liu Q, Fu Q M, et al. Batch least-squares policy iteration[J]. Computer Science, 2014, 41(9): 232-238.)
- [4] Ormoneit Dirk, Sen Saunak. Kernel-based reinforcement learning[J]. Machine Learning, 2002, 49(2): 161-178.
- [5] Xu Xin, Hu Dewen, Lu Xicheng. Kernel-based least squares policy iteration for reinforcement learning[J]. IEEE Trans on Neural Networks, 2007, 18(4): 973-992.
- [6] Xu Xin, Peng Chengzhang, Dai Bin, et al. A kernel-based reinforcement learning approach to stochastic pole balancing control system[C]. IEEE/ASME Int Conf on Advanced Intelligent Mechatronics. Montreal: IEEE Press, 2010: 1329-1334.
- [7] 程玉虎, 冯涣婷, 王雪松. 基于状态-动作图测地高斯基的策略迭代强化学习[J]. 自动化学报, 2011, 37(1): 44-51. (Cheng Y H, Feng H T, Wang X S. Policy iteration reinforcement learning based on geodesic gaussian basis defined on state-action graph[J]. Acta Automatica Sinica, 2011, 37(1): 44-51.)
- [8] Kim Min-soeng, Hong Gun-gi, Lee Ju-jang. On-line fuzzy Q-learning with extended rule interpolation technique[C]. IEEE/RSJ Int Conf on Intelligent Robots and Systems. Kyongju: IEEE Press, 1999, 2: 757-762.
- [9] Sun Haojun, Wang Shengrui, Jiang Qingshan. FCM-Based model selection algorithms for determining the number of clusters[J]. Pattern Recognition, 2004, 37(10): 2027-2037.
- [10] Kaufman Leonard, Rousseeuw Peter J. Finding groups in data: An introduction to cluster analysis[M]. New York: John Wiley and Sons Ltd, 1990: 68-125.
- [11] 于剑, 程乾生. 模糊聚类方法中的最佳聚类数的搜索范围[J]. 中国科学: E辑, 2002, 32(2): 274-280. (Yu J, Cheng Q S. The upper bound of the optimal number of clusters in fuzzy clustering[J]. Science in China: Series E, 2002, 32(2): 274-280.)

(责任编辑: 郑晓蕾)