

基于模糊聚类的绿色工艺评价样本分类方法

王宇钢¹, 修世超¹, 王柯元²

(1. 东北大学 机械工程与自动化学院, 辽宁 沈阳 110819; 2. 大连理工大学 电子信息与电气工程学部, 辽宁 大连 116024)

摘 要: 针对绿色工艺评价样本具有不确定性、多维性以及量纲差异大的特点,为实现样本的合理分类,提出一种基于核的模糊可能性聚类新算法.该方法将核模糊聚类算法、可能性聚类算法和减法聚类算法相结合,以提高聚类的准确率;使用聚类有效性指标作为分类条件,自适应确定最佳分类数.仿真实验结果表明,该算法具有较好的有效性和鲁棒性,并将该算法运用在绿色工艺评价样本分类中,得到了较好的分类效果,验证了算法的实用性.

关 键 词: 核模糊聚类;可能性聚类;减法聚类;有效性指标;绿色工艺;样本分类

中图分类号: TP 391 **文献标志码:** A **文章编号:** 1005-3026(2016)03-0387-05

Sample Classification Method for Green Process Evaluation Based on Fuzzy Clustering

WANG Yu-gang¹, XIU Shi-chao¹, WANG Ke-yuan²

(1. School of Mechanical Engineering & Automation, Northeastern University, Shenyang 110819, China; 2. Faculty of Electronic Information and Electrical Engineering, Dalian University of Technology, Dalian 116024, China. Corresponding author: WANG Yu-gang, E-mail: 9932783@qq.com)

Abstract: Due to the uncertainty, multidimensionality and significant difference of the evaluation samples of green process, a novel algorithm of kernel-based fuzzy possibilistic clustering was proposed in order to achieve reasonable sample classification. Kernel fuzzy clustering, possibilistic clustering and subtraction clustering were combined to improve the accuracy of clustering and cluster validity index was used as the classification condition to obtain the optimal classification number. The simulation results showed that this algorithm has good validity and robustness. When the algorithm is applied to classify the evaluation samples of green process, good classification effects are gained, which verifies its practicability.

Key words: kernel fuzzy clustering; possibilistic clustering; subtraction clustering; validity index; green process; sample classification

随着资源环境问题的日益严峻,对制造过程绿色特性进行评价,择优确定工艺方案已成为实现绿色制造的一种重要方法.国内外学者对绿色工艺评价方法进行了大量研究,提出了一些有效的评估方法,如生命周期评价法^[1]、模糊评判法^[2]、层次分析法^[3]等.但实际应用中,这些方法易出现评价过程繁琐、周期较长、过分依赖评价者主观判断等问题,从而导致评价结果的有效性和实用性变差.

机器学习方法具有很强的分析计算能力和泛

化能力,适合处理具有复杂性和模糊性特点的系统评价问题^[4],但采用机器学习方法进行绿色工艺评价的研究却鲜见报道.这主要由于要保证基于机器学习方法进行评价的效果,就需要有高质量的训练样本集.而工艺样本数据具有不确定性、多维性及量纲差异大的特点,使得训练样本的分类成为一项极其复杂的工作.

根据上述分析,提出一种基于核的模糊可能性聚类新算法(NKPFM),应用该算法对绿色工艺评价样本实现最优分类,最终为采用机器学习

收稿日期: 2015-01-07

基金项目: 国家自然科学基金资助项目(51375083).

作者简介: 王宇钢(1977-),男,辽宁锦州人,东北大学博士研究生;修世超(1958-),男,辽宁凌源人,东北大学教授,博士生导师.

方法进行绿色工艺评价提供决策支持. 该算法将核模糊聚类算法、可能性聚类算法及减法聚类算法相结合, 使用聚类有效性指标作为分类条件, 可对样本实现自适应分类. 仿真实验结果表明该算法可以实现自适应确定聚类数, 且具有较好的有效性和鲁棒性. 将该算法应用于绿色磨削工艺评价样本的分类, 取得较好效果.

1 模糊聚类算法

模糊聚类作为数据分析和建模的主要方法已得到广泛的应用, 其中主要有模糊 C-均值(FCM)聚类算法、可能性模糊 C-均值(PCM)聚类算法^[5]、基于核的模糊 C-均值(KFCM)聚类算法^[6]等. 但这些算法依然存在聚类数需预先设定, 聚类性能依赖初始聚类中心的选取等缺陷. 为克服这些缺陷, 提出一种基于核的模糊可能性聚类新算法.

1.1 基于核的模糊可能性聚类算法

基于核的模糊可能性聚类算法(KPFCM)利用核函数将在输入空间中线性不可分的样本在高维特征空间线性可分, 再利用可能性聚类放宽对隶属度的约束, 构造新的目标函数, 实现样本在高维特征空间中聚类^[7-8].

设样本集 $X = \{x_1, \dots, x_n\} \in \mathbf{R}^p$, 通过一个非线性映射函数 Φ 把所有样本映射到高维特征空间 F 中, 得到 $\Phi(x_1), \Phi(x_2), \dots, \Phi(x_n)$, 最终聚类在特征空间 F 中进行. 由核函数定义可知, 在原空间的点积运算可以表示为高维样本空间中核函数 $K(x, y)$ 的运算, $K(x, y) = \Phi(x) \Phi(y)$, 且满足 Mercer 条件: 对称性和 Carchy-Schwarz 不等式. 在高维特征空间中, KPFCM 算法目标函数表达式为

$$J_m(U, V) = \sum_{i=1}^C \sum_{k=1}^n u_{ik}^m \|\Phi(x_k) - \Phi(v_i)\|^2 + \sum_{i=1}^C \eta_i \sum_{k=1}^n (1 - u_{ik})^m. \quad (1)$$

式中: U 为模糊隶属度矩阵; V 为聚类中心; C 为聚类个数; x_k 为第 k 个样本; v_i 为第 i 个聚类中心; u_{ik} 为第 k 个样本属于第 i 类的隶属度; $m > 1$ 为加权指数; $\|\Phi(x_k) - \Phi(v_i)\|^2$ 表示核空间中第 j 个样本到第 i 个聚类中心的距离, 可表示为 $\|\Phi(x_k) - \Phi(v_i)\|^2 = K(x_k, x_k) + K(v_i, v_i) - 2K(x_k, v_i)$; η_i 是一个适合的正数, 可用式(3)计算:

$$\eta_i = K \sum_{k=1}^n 2u_{ik}^m (1 - K(x_k, v_i)) / \sum_{k=1}^n u_{ik}^m. \quad (3)$$

常见的核函数中, 高斯核函数由于对噪声点敏感度低而得到广泛应用. 高斯核函数表达式为

$$K(x, y) = \exp(-\|x - y\|^2 / \sigma^2). \quad (4)$$

式中, σ^2 为常数, 且 $K(x, x) = 1$. 此时, 目标函数表达式(1)可改写为

$$J_m(U, V) = 2 \sum_{i=1}^C \sum_{k=1}^n u_{ik}^m (1 - K(x_k, v_i)) + \sum_{i=1}^C \eta_i \sum_{k=1}^n (1 - u_{ik})^m. \quad (5)$$

式(5)达到最小值的条件为

$$v_i = \sum_{k=1}^n u_{ik}^m K(x_k, v_i) x_k / \left(\sum_{k=1}^n u_{ik}^m K(x_k, v_i) \right), \quad (6)$$

$$u_{ik} = \frac{(1 - K(x_k, v_i)) + (1 - K(x_k, v_i)) / \eta_i}{\sum_{j=1}^C (1 - K(x_k, v_j)) + (1 - K(x_k, v_i)) / \eta_i}^{-1/(m-1)}. \quad (7)$$

1.2 减法聚类

减法聚类使用数据样本密度函数计算数据点密度, 并把所有的数据点作为候选的聚类中心, 通过比较每个数据点的密度指标来确定该点作为聚类中心的可能性. 减法聚类过程如下:

1) 计算数据样本 (x_1, x_2, \dots, x_p) 中的每个数据点 x_i 的密度指标:

$$D_i = \sum_{j=1}^P \exp[-\|x_i - x_j\|^2 / (r_a/2)^2]. \quad (8)$$

式中, r_a 为聚类半径, 表示该点的一个邻域. 选择密度指标最大的数据点作为第一个聚类中心 D_{c1} .

2) 设 D_{ci} 为数据点 x_{ci} 的密度指标, 则每个数据点的密度指标按下式修正:

$$D_i = D_i - D_{ci} \exp(-\|x_i - x_{ci}\|^2 / (r_b/2)^2). \quad (9)$$

式中, r_b 为一个密度指标函数显著减小的邻域.

3) 在剩余的 $P-1$ 个数据点中选择密度指标最大的数据点作为新的聚类中心. 重复以上过程, 直到 D_{ci} 与 D_{c1} 的比值小于设定阈值, 聚类结束.

减法聚类的密度中心出现顺序依据密度值, 最早出现的聚类中心密度值最大, 且成为合适的聚类中心可能性最大. 因此若聚类个数为 C 时, 以减法聚类产生前 C 个聚类中心作为聚类的初始中心, 可有效避免算法对初始聚类中心敏感的问题.

1.3 聚类有效性指标

聚类有效性指标可用于确定样本划分的最佳聚类数目. KPFCM 算法与 FCM 算法一样需要预先指定聚类数目, 而在对数据集空间结构不了解

的情况下,预设的聚类数目很难保证为最合适的聚类数. 因此,通过应用合适的聚类有效性指标选择样本划分的聚类数,可实现最合理的聚类效果. Pakhira 等基于类内紧凑度和类间分离度定义的 PBMF 是近年来聚类性能较好的有效性指标^[9],其表达式如下:

$$\text{PBMF} = \frac{1}{C} \frac{E_1 \times \max_{i \neq j} \|v_i - v_j\|}{\sum_{i=1}^C \sum_{j=1}^n u_{ij}^m \|x_j - v_i\|} \quad (10)$$

式中: v_i 表示第 i 个聚类中心; x_j 表示第 j 个数据样本; E_1 为由数据集确定的常数. PBMF 指标值越大表示聚类效果越好.

1.4 改进后的算法

新的基于核的模糊可能聚类算法 (NKPFM) 具有两层迭代,内层迭代为减法聚类与 KPFM 算法的组合,通过最小化目标函数(5)实现聚类;外层迭代计算内层聚类结果的 PBMF 指标值,每迭代一次聚类数增加 1. NKPFM 算法过程如下:

- 步骤 1 设置参数:加权指数 m ,迭代终止阈值 ε ,高斯核函数参数 σ^2 ,由减法聚类获得最大聚类数 C_{\max} ,并令 $C=2$;
- 步骤 2 初始化聚类中心,并且依据公式(7)初始化隶属度矩阵;
- 步骤 3 根据公式(5)计算目标函数值,公式(6)更新聚类中心,公式(7)更新隶属度矩阵;
- 步骤 4 若相邻的目标函数值变化量达到阈值 ε ,则算法终止,否则转步骤 3.
- 步骤 5 根据式(10)计算 PBMF 值,若 $C < C_{\max}$,则 $C = C + 1$,转步骤 2;否则算法结束. 由 PBMF 最大值确定最佳聚类数和相应的样本分类结果.

2 仿真实验

分别运行 FCM,KFCM 和 NKPFM 算法,对 iris 和 wine 数据集进行测试. 其中,iris 数据集包含 4 维样本 150 个,分为三类. wine 数据集有 13 维样本 178 个,共分为三类,各含 59,71,48 个样本. 实验条件为:减法聚类半径 $r_a = r_b = 0.5$,误差 $\varepsilon = 0.000\ 01$,最大迭代次数 $T_{\max} = 100$, $m = 2.0$. 计算机配置:英特尔酷睿 2 双核 CPU,主频 2.20 GHz,内存 2.00 GB,利用 MATLAB 7.0 进行仿真实验.

2.1 聚类准确性

对 iris 数据集实验,运行 NKPFM 算法得到

聚类数为 3 时的隶属度函数分布如图 1 所示,相应的聚类中心为

$$V = \begin{bmatrix} 6.509\ 3 & 2.978\ 1 & 5.236\ 5 & 1.847\ 7 \\ 4.978\ 8 & 3.382\ 2 & 1.467\ 0 & 0.239\ 3 \\ 5.877\ 2 & 2.865\ 3 & 4.135\ 6 & 1.331\ 3 \end{bmatrix}.$$

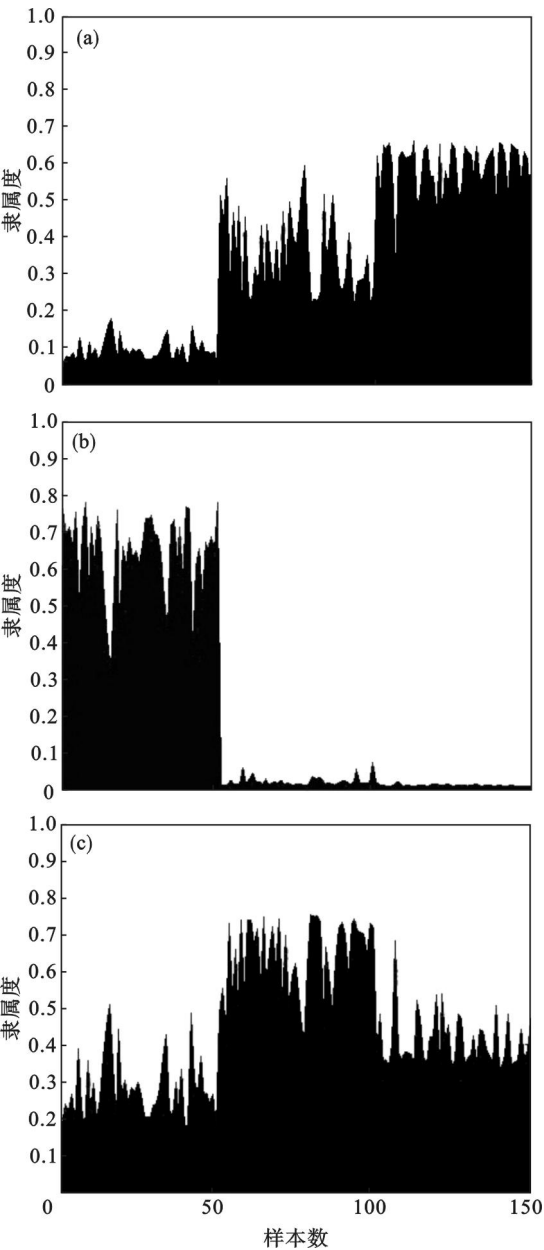


图 1 聚类后隶属度函数分布
Fig. 1 Membership function distribution after clustering
(a)一类一; (b)一类二; (c)一类三.

PBMF 指标值随聚类数的变化如表 1 所示. 当聚类数为 3 时,PBMF 有最大值,表明 iris 数据集的最佳聚类数为 3.

分别运行 FCM 算法、KFCM 算法和 NKPFM 算法 20 次,聚类结果如表 2 所示,括号外为样本典型值归类的误分数,括号内为模糊隶属度归类的误分数. 从表 2 可知,虽然 NKPFM

算法迭代次数和运行时间有所增加,但该算法误
分数明显低于另两种算法,聚类准确率更好.

表 1 Iris 数据集聚类 PBMF 的变化
Table 1 Change of PBMF for iris data set clustering

聚类数	2	3	4
PBMF	1. 287 8	1. 425 4	1. 193 8

表 2 Iris 数据集的聚类结果
Table 2 Clustering results of iris data set

算法	误分数	迭代次数	所需时间/s
FCM	20(26)	21	0. 094
KFCM	12(28)	11	0. 088
NKPFCM	8(15)	20	0. 098

为测试高维数据集聚类效果,对 wine 数据集进
行测试,测试结果如表 3 所示. 从表 3 可知,
NKPFCM 算法对高维数据集的聚类效果依然最好.

表 3 Wine 数据集的聚类结果
Table 3 Clustering results of wine data set

算法	误分数	迭代次数	所需时间/s
FCM	26(53)	51	0. 156
KFCM	47(70)	59	0. 172
NKPFCM	20(43)	51	0. 168

2.2 对噪声数据的处理

为考察样本的抗噪能力,分别对原样本集加
入 10%,20%,30%,40% 的噪声样本后,对比
KFCM 算法与 NKPFCM 算法的抗噪能力,结果如
图 2 所示.

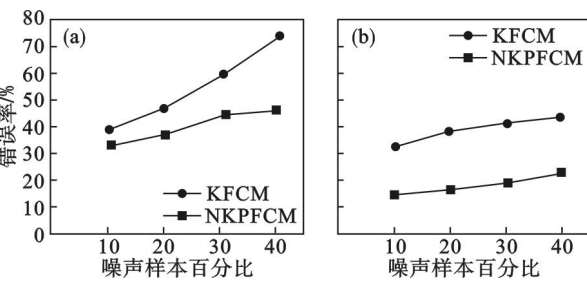


图 2 两种样本集的抗噪图
Fig. 2 Anti-noise-proof features of two sample sets
(a)—iris 数据集;(b)—wine 数据集.

由图 2 分析可知,两种算法的错误率相比加
入噪声数据前都有所增高,但是 NKPFCM 算法错
误率比 KFCM 算法的增长幅度小很多,且随着噪
声数据的增加,NKPFCM 算法错误增长率相比
KFCM 算法无大幅增加. 因此,NKPFCM 算法相
比 KFCM 算法具有更好的鲁棒性.

3 绿色工艺评价样本集的划分

以某汽车制造企业磨削加工工艺为例,采用
NKPFCM 算法对绿色磨削工艺评价样本进行分
类.

3.1 确定评价样本

本文选取与磨削工艺绿色度密切相关的 12
个指标作为样本评价指标,如表 4 所示. 由于评
价指标具有模糊性及量纲差异大的特点,为避免
专家主观判断的影响,本文采用定量分析和定性
描述相结合的半定量评价方法进行量化. 根据对
指标影响状况描述,将评价指标分为定性指标和
定量指标,对于无法计量的定性指标,由专家采
用十分制打分进行量化;对于定量指标则直接
采用测量值. 具体的指标描述和评分方法可参
见文献 [10].

表 4 磨削工艺评价指标
Table 4 Grinding process evaluation indexes

定量指标	定性指标
X_1 磨削液毒性、 X_2 刺激 气味、 X_4 粉尘、 X_7 夹具 类型、 X_{12} 安全性	X_3 噪声、 X_5 原材料消耗、 X_6 磨具消耗、 X_8 磨削液用量、 X_9 能耗、 X_{10} 加工费、 X_{11} 折 旧费

磨削工艺评价样本量化后数据见表 5.

表 5 磨削工艺样本量化值
Table 5 Quantized values of grinding process samples

n	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}	X_{11}	X_{12}
1	3	5	78	3	0.018	4.5	1	6.1	5.4	42.0	1.4	1
2	7	5	80	5	0.012	3.0	3	5.2	5.1	40.0	1.2	3
3	1	3	75	3	0.009	3.0	1	5.0	4.8	35.0	1.1	1
4	5	7	82	5	0.016	4.0	3	5.5	5.6	45.0	1.4	3
5	5	5	78	5	0.014	3.5	3	5.2	5.2	39.0	1.3	1
6	5	5	85	3	0.015	3.5	3	6.0	5.9	42.5	1.2	1
7	3	7	80	3	0.019	4.5	1	5.4	5.5	40.5	1.4	1
8	7	5	78	7	0.011	3.0	3	5.0	5.0	40.0	1.1	3
9	7	5	80	5	0.016	4.0	3	5.5	5.3	37.0	1.3	3
10	3	7	65	3	0.018	4.5	1	6.0	6.2	35.5	1.4	1
11	7	5	86	7	0.013	4.5	3	5.4	5.3	40.5	1.2	3
12	5	5	68	5	0.015	3.5	3	5.8	6.0	37.5	1.3	1
13	3	5	65	5	0.017	4.0	3	6.0	5.4	35.0	1.4	1
14	3	7	78	5	0.020	5.0	1	6.5	6.0	50.0	1.5	1
15	7	5	90	5	0.011	3.0	3	5.5	5.2	39.0	1.1	3
16	3	5	65	5	0.019	4.5	1	6.5	6.3	38.5	1.4	1
17	3	7	80	3	0.018	4.5	3	6.3	5.9	54.5	1.4	1
18	7	9	76	5	0.010	3.0	5	5.0	5.2	48.0	1.1	3
19	3	5	76	5	0.017	4.0	1	6.0	5.9	38.5	1.4	1
20	5	5	69	5	0.016	4.0	5	5.6	5.5	42.0	1.3	1

3.2 样本的划分

对量化数据采用式(11)进行归一化处理:

$$x'_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}. \tag{11}$$

其中: x_i 为指标测量值; x_{\min}, x_{\max} 分别为数据集中该指标最小值与最大值. 利用 NKPFM 算法对数据进行聚类,PBMF 值的变化如表 6 所示,并依据 PBMF 最大值确定聚类数目为 3.

表 6 数据集聚类 PBMF 的变化

Table 6 Change of PBMF for data set clustering					
聚类数	2	3	4	5	6
PBMF	0.385 9	0.974 4	0.565 6	0.543 5	0.487 62

聚类数设为 3, KPFCM 算法和 NKPFM 算法的样本聚类结果如表 7 所示. 表 7 显示样本 3, 7, 19 相对两种算法的不同分类结果, 比较两种算法分类性能, NKPFM 算法的 PBMF 值较大, 表明 NKPFM 算法的聚类结果更优秀, 如表 8 所示. 因此, 采用 NKPFM 算法获得的分类样本可以作为训练样本集为基于机器学习方法进行绿色工艺评价提供决策支持.

表 7 样本聚类结果
Table 7 Sample clustering results

算法	样本号
KFCM	类一: 3, 10, 12, 13, 16, 20
	类二: 2, 5, 6, 7, 9, 11, 15, 19
	类三: 1, 4, 8, 14, 17, 18
NKPFM	类一: 10, 12, 13, 16, 20
	类二: 2, 5, 6, 9, 11, 15
	类三: 1, 3, 4, 7, 8, 14, 17, 18, 19

表 8 分类性能比较

Table 8 Classification performance comparison			
算法	迭代次数	目标函数误差	PBMF
KFCM	28	0.000 008	0.905 721
NKPFM	31	0.000 003	0.974 402

4 结 论

1) 构造了一个新的基于核的模糊可能性聚类算法(NKPFM), 经过对 iris 和 wine 数据集的仿真测试, 结果表明该算法具有较好的准确性和

鲁棒性.

2) 阐述了磨削加工的绿色工艺评价样本集的生成方法, 并将 NKPFM 算法应用于绿色工艺评价样本的划分, 获得了较好的分类效果.

3) 有效的分类样本可以对采用机器学习方法的绿色工艺评价提供决策支持, 接下来还需进一步研究如何选择合适的机器学习评价方法.

参考文献:

[1] Remo A P, Diogo A L, Eraldo J, et al. Dynamic system for life cycle inventory and impact assessment of manufacturing processes[C]//The 21st CIRP Conference on Life Cycle Engineering. Trondheim: Procedia CIRP, 2014: 531 – 536.

[2] 王桂萍, 贾亚洲, 周广文. 基于模糊可拓层次分析法的数控机床绿色度评价方法及应用[J]. 机械工程学报, 2010, 46(3): 141 – 147.

(Wang Gui-ping, Jia Ya-zhou, Zhou Guang-wen. Evaluation method and application of CNC machine tool's green degree based on fuzzy-EAHP [J]. Journal of Mechanical Engineering, 2010, 46(3): 141 – 147.)

[3] Ng C Y, Chuah K B. Evaluation of design alternatives' environmental performance using AHP and ER approaches [J]. IEEE Systems Journal, 2013, 8(4): 1182 – 1189.

[4] Qiao L, Rajagopalan C, Clifford G D. Ventricular fibrillation and tachycardia classification using a machine learning approach[J]. Biomedical Engineering, 2014, 61(6): 1617 – 1603.

[5] Yang H D, Li C S, Hu J. RFID intrusion detection with possibilistic fuzzy c-means clustering [J]. Journal of Computational Information Systems, 2010, 6(8): 2623 – 2632.

[6] Sadaaki M. Different objective functions in fuzzy c-means algorithms and kernel-based clustering [J]. International Journal of Fuzzy Systems, 2011, 13(2): 89 – 97.

[7] Tushir M, Srivastava S. A new kernel based hybrid c-means clustering model [C]//Proceedings of 2007 IEEE International Conference on Fuzzy Systems. London: IEEE, 2007: 1 – 5.

[8] Lucieer V, Lucieer A. Fuzzy clustering for seafloor classification[J]. Marine Geology, 2009, 264(3/4): 230 – 241.

[9] Pakhira M K, Bandyopadhyay S, Maulik U. Validity index for crisp and fuzzy clusters[J]. Pattern Recognition, 2004, 37(3): 487 – 501.

[10] 刘飞. 绿色制造的理论与技术[M]. 北京: 科学出版社, 2005: 144 – 151.

(Liu Fei. Theory and technology of green manufacturing [M]. Beijing: Science Press, 2005: 144 – 151.)