

# 涡轮泵试车数据单类支持向量机检测算法<sup>\*</sup>

胡 雷, 胡 莺 庆, 秦 国 军

(国防科技大学 机电工程与自动化学院, 湖南 长沙 410073)

**摘要:** 为了在缺乏故障样本的情况下检测某型液体火箭发动机涡轮泵故障, 实现基于不完整信息的状态决策, 建立了基于  $v$ -支持向量分类器的单类支持向量机新异类检测模型。在分析了模型决策边界、支持向量和约束条件之间关系的基础上, 为单类支持向量机引入并改进了序贯最小优化算法, 提高了训练效率, 解决了大样本训练问题。通过对某型液体火箭发动机涡轮泵历史试车数据的分析, 结果表明, 所建模型的训练速度得到了很大提高, 对涡轮泵状态的检测效果良好。

**关键词:** 液体推进剂火箭发动机; 涡轮泵; 新异类检测模型<sup>+</sup>; 单类支持向量机<sup>+</sup>; 序贯最小优化<sup>+</sup>; 故障检测

中图分类号: V434.21 文献标识码: A 文章编号: 1001-4055 (2008) 02-0244-05

## Support vector machines detection method for turbopump test data analysis

HU Lei HU Niao-qing QIN Guo-jun

(Inst. of Mechatronics Engineering and Automation, National Univ. of Defense Technology, Changsha 410073, China)

**Abstract** For lacking of fault samples, it is very difficult to detect the faults of a Liquid Rocket Engine (LRE) turbopump and make decision based on incomplete information. To solve this problem, a  $v$ -support vector machine novelty detection model was founded. Taking into account of the relationship between decision boundary, support vectors and constraints, a training algorithm based on Sequential Minimal Optimization (SMO) was introduced and improved for One-Class Support Vector Machines (OCSVM). With the analysis of LRE historical test data, it showed that SMO algorithm improves the training efficiency evidently and enables the model to deal with large training data. And this model trained by SMO can detect the faults of the LRE turbopump well.

**Key words** Liquid propellant rocket engine; Turbine pump; Novelty detection model; One-class support vector machines<sup>+</sup>; Sequential minimal optimization<sup>+</sup>; Fault detection

## 1 引言

液体火箭发动机涡轮泵工作在极端的转子动力条件下, 是液体火箭发动机的故障多发部位<sup>[1,2]</sup>。目前, 针对涡轮泵故障检测算法的研究很多, 但这些算法多数都需要足够多的正常数据和异常数据进行训练。由于故障数据的获取往往意味着巨大的损失, 所以实际的试车数据多数都是正常数据, 故障数据相对

稀少。因此, 为涡轮泵试车数据分析建立无需故障数据参与训练的新异类检测 ND (Novelty Detection) 模型是十分必要的。

ND 技术中的单类支持向量机 OCSVM (One-Class Support Vector Machines) 方法是支持向量机 SVM (Support Vector Machines) 的变种, 和 SVM 一样, 它完美地结合了结构风险最小化、最优分类超平面、核到内积空间中的映射、不可分问题中的松弛变量等

\* 收稿日期: 2007-02-15 修订日期: 2007-06-18。

基金项目: 国家自然科学基金 (50375153); 高等学校全国优秀博士学位论文专项资金 (200434)。

作者简介: 胡雷 (1981—), 男, 博士生, 研究领域为机器状态监控与故障诊断、机械信号处理等。

E-mail lake\_h@mail.com

思想, 具有传统模式识别方法(包括神经网络)所不具备的许多优良性能, 如无局部最小值、能够解决非线性学习问题、推广能力强等<sup>[3,4]</sup>。

前期为涡轮泵试车数据分析建立了基于支持向量数据描述的 OCSVM 检测模型<sup>[5]</sup>。但是, 同 SVM一样, OCSVM 的训练也需求解一个二次优化问题, 而该二次优化问题会因需要计算和存储核函数矩阵而变得缓慢甚至无法进行。为解决这一问题, 本文引入 Platt J 为 SVM 设计的序贯最小优化 SMO ( Sequential M inimal Optimization)<sup>[6]</sup>思想, 提出了用以解决 OCSVM 训练问题的 SMO 算法, 并研究了其中的初始值设置方法, 以及模型循环和更新方法。

## 2 OCSVM 检测模型

这里使用基于  $v$ -支持向量分类器  $v$ -SVC ( $v$ -Support Vector Machine)<sup>[7]</sup> 的 ND 模型。 $v$ -SVC 是一种 OCSVM 超平面模型, 其思想是首先使用核函数映射  $\Phi$ , 将训练样本  $x$ (又称目标类样本)投影到一个高维特征空间, 然后在特征空间中构造分类超平面  $f(x)$ , 将训练样本与原点以最大间隔分开, 同时要求尽可能少的训练样本位于原点一侧, 如图 1 所示。

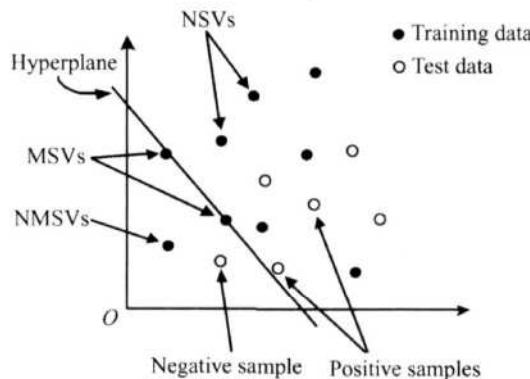


Fig 1 OCSVM model based on  $v$ -SVC

$v$ -SVC 的训练模型为

$$\text{minimize } P(\alpha) = \sum_{i,j=1}^l \alpha_i \alpha_j K(x_i, x_j) \quad (1)$$

$$\text{subject to } 0 \leq \alpha_i \leq 1/vl \quad i = 1, \dots, l \quad (2)$$

$$\sum_{i=1}^l \alpha_i = 1 \quad (3)$$

该模型是在约束条件(2)和(3)下, 求使得目标函数  $P(\alpha)$  取最小值的拉格朗日乘子  $\alpha$ 。其中  $x_i$  和  $x_j$  是训练样本,  $l$  是样本个数,  $v$  是用户设置的平滑参数, 用以平衡训练误差和推广能力,  $K$  是核函数, 这里使用高斯核函数

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{S^2}\right) \quad (4)$$

式中  $S$  是高斯参数。

可以证明平滑参数  $v$  是目标类上错误率的上界, 是支持向量占训练样本百分率的下界<sup>[8]</sup>。文献[9]讨论了参数  $v$  对 OCSVM 影响, 文献[10]分析了  $v$  的变化给涡轮泵试车数据分析结果造成的影响, 并给出了确定  $v$  和高斯参数  $s$  的一般方法。

分别称满足  $f(x) > 0$ ,  $f(x) = 0$  和  $f(x) < 0$  的训练样本为非支持向量 NSVs、边界支持向量 MSVs 和非边界支持向量 NMSVs(如图 1 所示), 对应的拉格朗日乘子分别满足  $\alpha_i = 0$ ,  $0 < \alpha_i < 1/vl$  和  $\alpha_i = 1/vl$ 。检测时将满足  $f(x) > 0$  的测试样本判为正常样本, 将满足  $f(x) < 0$  的样本判为异类样本。模型的测试函数为

$$f(x) = \sum_k \alpha_k K(x_k, x) - b \quad (5)$$

其中  $x$  是测试数据,  $x_k$  属于 MSVs 或 NMSVs,  $\alpha_k$  为与之对应的拉格朗日乘子。由于 MSVs 满足  $f(x) = 0$ , 所以任选  $x_m \in \text{MSVs}$  可得阈值

$$b = \sum_k \alpha_k K(x_k, x_m) \quad (6)$$

## 3 OCSVM 的 SMO 算法

SMO 算法的每一步迭代仅选择两个数据点对应的拉格朗日乘子进行优化, 其优势在于两个数据点的优化可以获得解析解, 从而不需将二次优化规划看作算法的一部分<sup>[3]</sup>。本文为 OCSVM 设计的 SMO 算法如下。

### 3.1 $\alpha_1$ 和 $\alpha_2$ 的优化解析解

选定每步要更新的两个拉格朗日乘子  $\alpha_1$  和  $\alpha_2$ , 其余  $\alpha_i$  ( $i = 3, \dots, l$ ) 不变, 则目标函数为

$$P(\alpha_1, \alpha_2) = \alpha_1^2 K_{11} + \alpha_2^2 K_{22} + 2\alpha_1 \alpha_2 K_{12} + 2 \sum_{i=3}^l \alpha_i \alpha_i K_{1i} + 2 \sum_{i=3}^l \alpha_i \alpha_i K_{2i} + \text{const} \quad (7)$$

式中  $K_{ij}$  表示  $K(x_i, x_j)$ 。由约束条件(3)知,  $\alpha_1 + \alpha_2$  不变。这里令

$$r = \alpha_1 + \alpha_2 = \alpha_1^{old} + \alpha_2^{old} \quad (8)$$

$$U = [U_j]_{l \times 1} = [\sum_{i=1}^l \alpha_i K_{ij}]_{l \times 1} \quad (9)$$

$$V_j = \sum_{i=3}^l \alpha_i^{old} K_{ij} = U_j^{old} - \alpha_1^{old} K_{1j} - \alpha_2^{old} K_{2j} \quad (10)$$

式中  $U_j^{old} = \sum_{i=1}^l \alpha_i^{old} K_{ij}$ 。将式(10)代入式(7)得

$$P = P(\alpha_1, \alpha_2) = \alpha_1^2 K_{11} + \alpha_2^2 K_{22} + 2\alpha_1 \alpha_2 K_{12} + 2V_1 \alpha_1 + 2V_2 \alpha_2 + \text{const} =$$

$$\begin{aligned}
 & (r - \alpha_2)^2 K_{11} + \alpha_2^2 K_{22} + 2(r - \alpha_2) \alpha_2 K_{12} + \\
 & 2V_1(r - \alpha_2) + 2V_2 \alpha_2 + \text{const} = \\
 & (K_{11} + K_{22} - 2K_{12}) \alpha_2^2 + (2K_{12} - 2K_{11} + \\
 & 2V_2 - 2V_1) \alpha_2 + \text{const} \quad (11)
 \end{aligned}$$

$P$  在  $\alpha_2$  处取最小值, 应满足  $\partial P / \partial \alpha_2 = 0$ ,  $\partial^2 P / \partial \alpha_2^2 > 0$  由  $\partial P / \partial \alpha_2$  可得

$$\begin{aligned}
 \alpha_2^{\text{new, unc}} &= (\kappa K_{12} - \kappa K_{11} + V_2 - V_1) / (2K_{12} - \\
 K_{11} - K_{22}) = \alpha_2^{\text{old}} + (U_2^{\text{old}} - U_1^{\text{old}}) / (2K_{12} - \\
 K_{11} - K_{22}) \quad (12)
 \end{aligned}$$

而由  $\partial^2 P / \partial \alpha_2^2 > 0$  可得  $2K_{12} - K_{11} - K_{22} < 0$  该条件只要  $x_1 \neq x_2$  就能满足。

根据式(2), (3)和(8)知  $\alpha_2^{\text{new}}$  需满足

$$\begin{aligned}
 \max(0, \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C) \leq \alpha_2^{\text{new}} \leq \\
 \min(C, \alpha_1^{\text{old}} + \alpha_2^{\text{old}}) \quad (13)
 \end{aligned}$$

所以要对  $\alpha_2^{\text{new, unc}}$  进行剪辑

$$\alpha_2^{\text{new}} = \begin{cases} L & \text{若 } \alpha_2^{\text{new, unc}} < L \\ \alpha_2^{\text{new, unc}} & \text{若 } L < \alpha_2^{\text{new, unc}} < H \\ H & \text{若 } \alpha_2^{\text{new, unc}} > H \end{cases} \quad (14)$$

其中

$$\begin{aligned}
 L &= \max(0, \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C) \\
 H &= \min(C, \alpha_1^{\text{old}} + \alpha_2^{\text{old}})
 \end{aligned}$$

而由  $\alpha_1 + \alpha_2 = \alpha_1^{\text{old}} + \alpha_2^{\text{old}}$  可得

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - \alpha_2^{\text{new}} \quad (15)$$

式(14)和(15)就是所选两数据点对应的拉格朗日乘子的更新方程。

### 3.2 $\alpha_1$ 和 $\alpha_2$ 的循环选择与停机准则

使用启发式的选择算法选择  $\alpha_1$  和  $\alpha_2$ 。算法分两层循环: 外循环分别在所有样本和 MSVs 中交替进行, 搜索违反 KKT 条件的样本点, 即不满足条件

$$f(x) = \begin{cases} > 1 & \text{当 } \alpha_i = 0 \\ = 1 & \text{当 } 0 < \alpha_i < 1/vl \\ < 1 & \text{当 } \alpha_i = 1/vl \end{cases}$$

的样本点对应的拉格朗日乘子作为  $\alpha_2$ 。内循环寻找使得目标函数  $P$  发生最大变化的点对应的拉格朗日乘子作为  $\alpha_1$ , 一个快速的启发式方法是选择使得  $|U_2^{\text{old}} - U_1^{\text{old}}|$  最大的  $\alpha_1$ 。如果这个选择没有使得目标函数产生大的增长, 则从随机位置开始轮流尝试每个 MSVs, 如果仍不能引起目标函数产生大的增长, 则从随机位置开始轮流尝试所有样本。

循环过程中, 当所有样本都满足 KKT 条件时, 循环停止。为保证算法快速收敛, 选择较低的容忍等

级, 这样并不会给模型精度造成大的影响<sup>[3]</sup>。

### 3.3 循环过程中的更新

循环过程中, 除了要按照式(14)和(15)对拉格朗日乘子进行更新外, 还需要更新变量  $U$ ,  $b$ 。因为在选择待优化的乘子和判断停机准则时, 需要使用  $U$  和  $b$ 。

记  $\Delta \alpha_1 = \alpha_1^{\text{new}} - \alpha_1^{\text{old}}$ , 根据式(15), 有  $\Delta \alpha_1 = \alpha_1^{\text{new}} - \alpha_1^{\text{old}} = -\Delta \alpha_2$ 。再由式(9)可得  $U$  的更新方程

$$U_j^{\text{new}} = \sum_{i=1}^l \alpha_i^{\text{old}} K_{ij} + \Delta \alpha_2 K_{2j} + \Delta \alpha_1 K_{1j} = U_j^{\text{old}} + \Delta \alpha_2 K_{2j} - \Delta \alpha_1 K_{1j} \quad (16)$$

根据式(6), 若  $\alpha_2^{\text{new}}$  满足  $0 < \alpha_2^{\text{new}} < 1/vl$  则  $x_2 \in \text{MSVs}$  有

$$b = \sum_k \alpha_k K_{k2} \quad (17a)$$

否则, 考查  $\alpha_1^{\text{new}}$  是否满足  $0 < \alpha_1^{\text{new}} < 1/vl$  若满足, 则  $x_1 \in \text{MSVs}$  有

$$b = \sum_k \alpha_k K_{k1} \quad (17b)$$

而若  $x_1, x_2$  都不属于 MSVs 则

$$b = \text{mean}(\sum_k \alpha_k K_{km}) \quad (17c)$$

式(17)中,  $0 < \alpha_k \leq 1/vl$   $x_m$  是 MSVs

Keerthi 提出了一种改进的 SMO 方法<sup>[11]</sup>, 该方法在判断 KKT 条件时不需要使用  $b$  也就不需对  $b$  进行更新。事实上, 更新  $U$  和  $b$  时, 计算量主要体现在计算  $K_{2j}$  和  $K_{1j}$  上, 而更新  $U$  时必须计算  $K_{2j}$  和  $K_{1j}$ , 在计算了  $K_{2j}$  和  $K_{1j}$  的基础上再更新  $b$  对计算量影响不大。

### 3.4 初始值的设置

Platt J 为 SVM 设计的 SMO 算法中将所有拉格朗日乘子集的初始值设置为零。但是对于 OCSVM, 拉格朗日乘子集的初始值必须满足式(3)。考虑到决定模型边界的是支持向量, 且这些支持向量多位于训练样本分布区域的外围, 所以选择位于训练样本分布区域外围的点为初始支持向量。具体地, 分别选择每一维特征的最大值和最小值对应的样本点, 假设样本维数为  $n$  则选择了  $2n$  个样本(若某一特征的一个最值对应多个样本, 则选择其中的一个), 去掉其中重复的点, 得到  $m$  个样本, 可令这  $m$  个样本为初始的支持向量, 初始化其对应的拉格朗日乘子为  $1/m$ , 其余点对应的拉格朗日乘子为 0。

初始化拉格朗日乘子之后,  $U$  和  $b$  的初始值可以分别由式(9)和式(6)计算。但是初始选择的支持向量并不一定满足  $f(x) = 0$  所以取  $b = \text{mean}_k(\sum_i \alpha_i K_{ik})$ , 其中  $x_k$  就是初始选择的支持向量。

## 4 涡轮泵试车数据分析

### 4.1 涡轮泵试车数据检测

液体火箭发动机涡轮泵转速高, 载荷大, 振动信号中伴有大量的环境噪声干扰, 因此信号时域特征参数的一致性较差。相比之下, 频域特征特别是功率谱特征的一致性较好, 且能在一定程度上突出周期频率成份而抑制随机振动频率成份。因此提取振动信号的功率谱特征进行检测。

选取氢涡轮泵的径向振动信号进行分析。在振动信号序列上, 每 8 192 点(历时约 0.16 s, 信号的采样频率为 50 kHz)进行一次功率谱估计。对多次试车振动信号的功率谱分析结果表明, 信号的功率谱成份主要集中分布在 0~1 500 Hz 之间。将 0~1 500 Hz 之间的功率谱成分划分为 30 个频率段, 求各段的谱面积与总体谱面积之比, 然后作归一化处理, 构成 30 维的频段能力比特征向量。图 2 和图 3 所示分别为 T01 次试车故障前后的频段能量比  $r$  的分布情况。

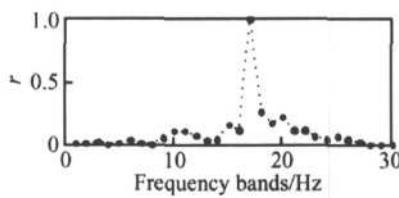


Fig. 2 Frequency band energy ratio  $r$  of T01 in normal condition

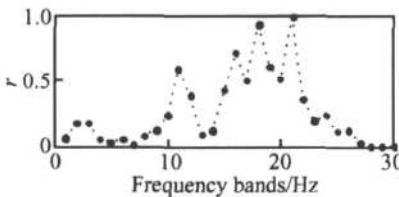


Fig. 3 Frequency band energy ratio  $r$  of T01 in faulty condition

由于液体火箭发动机推力室和燃气发生器内的压力脉动以及涡轮泵的高速旋转, 使发动机始终工作在强烈的振动状态中, 且涡轮泵还工作在高温、高压、高能量集中的恶劣环境中, 因此涡轮泵故障多是突发性故障<sup>[2]</sup>。因此可假设每次试车稳态过程的起始段是正常的。设置高斯参数  $s=20$ , 平滑参数  $v=0.05$ , 用稳态过程前 8 s 振动信号的频段能量比特征(50 个)作为训练样本, 训练基于 SMO 的  $v$ -SVC 模型, 得到一组支持向量和支持向量所对应的拉格朗日乘子  $\alpha_k$  以及阈值  $b$ , 使用该训练结果对剩余时间振动信号

的频段能量比特征进行检测。图 4 和图 5 所示分别为 T01 次试车和 T02 次试车振动信号的检测结果  $f(x_i)$ , 大于 0 时为正常, 否则为故障。

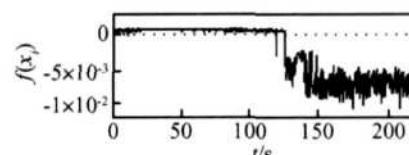


Fig. 4 Detection results of T01

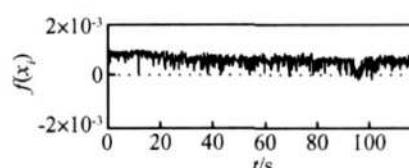


Fig. 5 Detection results of T02

### 4.2 检测效率分析

使用所建  $v$ -SVC 单类支持向量机模型对 5 次试车的氢泵径向振动信号进行检测。选择相同的模型参数:  $s=20$ ,  $v=0.05$ , 分别使用 SMO 算法和二次优化算法(采用 Matlab 的 quadprog 模块)对 5 次试车的振动信号特征进行优化训练和检测, 检测误差如表 1 所示。表中

$$F_+ = \frac{\text{异常的正常样本个数}}{\text{正常样本总数}}$$

$$F_- = \frac{\text{正常的异常样本个数}}{\text{异常样本总数}}$$

从表中可以发现, 针对涡轮泵振动信号, 使用 SMO 算法的检测精度更高。当然, 如果对使用二次优化算法训练得到的模型的阈值进行适当的缩放, 可以得到更好的检测精度<sup>[5]</sup>。

使用两种优化方法分别对不同长度的样本各进行 50 次训练, 耗时如表 2 所示。从表中可以看出, 与二次优化算法相比, SMO 算法的计算速度有明显提高; 而且二次优化算法的耗时受样本长度的影响很大, 甚至会因样本过大而无法得到优化结果, 而 SMO 算法的耗时受样本长度影响非常小。导致这一结果的主要原因有两点: 一是二次优化算法需要计算和存储核函数矩阵, 而仅仅存储核矩阵就需要一个随样本大小二次增长的内存空间, 而 SMO 算法不需矩阵运算。二是虽然 SMO 算法使用了两层循环, 但是每次循环只需进行很少的运算, 特别是使用了 3.4 节中的初始值设置方法后, 算法收敛速度很快。

**Table 1 Detection error rate with quadprog module and SMO**

Test No	Actual condition	Optimization method	$F_+$ (False negative rate)	$F_-$ (False positive rate)
T01	Faulty after 125 s	quadprog	7.72%	0
		SMO	0.26%	0
T02	Normal	quadprog	15.56%	—
		SMO	2.64%	—
T03	Normal	quadprog	9.44%	—
		SMO	3.19%	—
T04	Normal	quadprog	0.56%	—
		SMO	0	—
T05	Normal	quadprog	0	—
		SMO	0	—

**Table 2 Time consumed using two optimization methods**

Number of samples	Optimization method	Longest time consumed/s	Mean time consumed/s
50	quadprog	0.4063	0.3531
	SMO	0.0313	0.0275
80	quadprog	0.9531	0.8425
	SMO	0.0469	0.0291
120	quadprog	2.8125	2.1812
	SMO	0.0579	0.0388
200	quadprog	11.8281	9.1294
	SMO	0.0781	0.0628

## 5 结 论

基于  $v$ -SVC 的单类支持向量机新异类检测模型无需异常样本参与训练便可以得到涡轮泵振动信号的决策区域, 而基于序贯最小优化的单类支持向量机优化算法可以显著提高训练效率, 解决大样本训练问题。对某型液体火箭发动机涡轮泵历史试车数据的分析结果表明, 基于 SMO 的单类支持向量检测算法可以高效地识别涡轮泵的正常状态和故障状态。

## 参考文献:

- [1] 谢光军, 胡茑庆, 胡雷. 涡轮泵实时故障检测的改进自适应相关阈值算法 [J]. 推进技术, 2006, 27(1). (XIE Guang-jun, HU Niao-qing, HU Lei. Improved adaptive correlation thresholds algorithm for turbopump real-time fault detection [J]. *Journal of Propulsion Technology*, 2006, 27(1). )
- [2] 谢光军. 液体火箭发动机涡轮泵实时故障检测技术及系统研究 [D]. 长沙: 国防科技大学, 2006
- [3] Nello Cristianini; John Shawe-Taylor著. 李国政译. 支持向量机导论 [M]. 北京: 电子工业出版社, 2004
- [4] Vapnik Vladimir N. 统计学习理论的本质 [M]. 张学工译. 北京: 清华大学出版社, 2000
- [5] 胡雷, 胡茑庆, 谢光军. 基于单类支持向量机的涡轮泵故障检测方法研究 [C]. 2006年全国振动工程及应用学术会议论文集, 2006
- [6] Platt J. Sequential minimal optimization: a fast algorithm for training support vector machines [R]. *Technical Report MSR-TR-98-14 Microsoft Research*, 1998
- [7] Bernhard Schölkopf, Robert Williamson, Alex Smola, et al. Support vector method for novelty detection [C]. *Advances in Neural Information Processing Systems 12: Proceedings of the 1999 Conference*, MIT Press, 2000, 582~588
- [8] Chu Shun-kwong. Scaling up support vector data description by using core-sets [D]. *Department of Computer Science, Hong Kong University of Science and Technology*, 2004
- [9] Tax D M J, Duin R P W. Support vector data description [J]. *Machine Learning*, 2004, 54(1).
- [10] 胡雷. 涡轮泵试车数据分析及新异类状态检测技术研究 [D]. 长沙: 国防科技大学, 2005
- [11] Keerthi S S, Shevade S K, Bhattacharyya C, et al. Improvements to platt's SMO algorithm for SVM classifier design [J]. *Neural Computation*, 2001, 13(3).

(编辑: 郭振伶)