

# 基于小球大间隔方法的机械故障检测

郝腾飞 陈 果

南京航空航天大学,南京,210016

**摘要:**针对机械故障检测中,正常样本多、故障样本少、训练样本严重不平衡的客观情况,将小球大间隔方法引入其中,提出了一种不平衡样本下的机械故障检测方法。该方法同时使用大量的正常样本和少量的故障样本进行训练,在特征空间中构造一个包围正常样本的超球,在该超球体积最小化的同时,进一步使超球边界与故障样本之间的间隔最大化,从而显著减小将故障情况误判为正常情况的概率。将该方法应用到滚动轴承故障检测中,并与支持向量机和支持向量数据描述方法进行了比较,实验结果表明,该方法在解决不平衡样本下机械故障检测问题具有优越性。

**关键词:**故障检测;不平衡样本;小球大间隔;支持向量机;支持向量数据描述

中图分类号:TH17; TP206.3

DOI:10.3969/j.issn.1004-132X.2012.15.001

## Machinery Fault Detection Based on a Small Sphere and Large Margin Approach

Hao Tengfei Chen Guo

Nanjing University of Aeronautics and Astronautics, Nanjing, 210016

**Abstract:** In machinery fault detection, normal examples are much more than fault examples and the training examples are highly imbalanced. Aiming at this problem, a small sphere and large margin approach was used for machinery fault detection and a machinery fault detection method for imbalanced examples was put forward. The proposed method can use both of many normal examples and few fault examples to train. It constructed a hypersphere that contained normal examples in the feature space by training, such that the volume of this sphere was as small as possible, while at the same time the margin between the surface of this sphere and the fault examples was as large as possible. This method was applied to fault detection of rolling element bearings and comparisons were conducted with support vector machine and support vector data description. Experimental results validate its effectiveness in the machinery fault detection where the training examples are highly imbalanced.

**Key words:** fault detection; imbalanced dataset; small sphere and large margin; support vector machine; support vector data description

## 0 引言

机械故障检测在本质上是一个模式识别问题,建立在统计学习理论之上的支持向量机(support vector machines, SVM)<sup>[1-2]</sup>具有良好的推广能力,已经在机械故障检测领域得到了广泛应用<sup>[3-8]</sup>。但是,支持向量机作为一种两类分类方法,在训练中须同时使用正常样本和故障样本。机械故障检测中,故障样本通常很难获取,也不可能为了获得故障样本而故意破坏机械设备,因此,在机械故障检测中,故障样本是可遇而不可求的。针对该问题,一些学者将支持向量数据描述(support vector data description, SVDD)<sup>[9-10]</sup>应用于机械故障检测<sup>[11-13]</sup>,该方法只需使用正常样本进行训练,因此有效地解决了故障样本缺失情况下的故障检测问题。但是,在现实的机械故障检测中,故障样本虽然不易获取,但一般通过各种途径还是能获取到一些,如通过机械设备曾经偶尔发

生的一些故障可以收集到一些故障样本,只是这些样本相对于正常样本较少,因此,机械故障检测的现实情况是正常样本较多,故障样本较少,两者在数量上严重不平衡。在这种情况下,如果使用传统的支持向量机进行故障检测,由于训练样本严重不平衡,其性能会显著下降。如果使用支持向量数据描述进行故障检测,则故障样本得不到有效利用。基于上述分析,在机械故障检测领域,研究不平衡样本下的故障检测方法是一个重要且有意义的问题。

针对该问题,本文将小球大间隔方法(small sphere large margin, SSLM)<sup>[14]</sup>应用于机械故障检测,提出了一种不平衡样本下的机械故障检测方法。该方法在训练中不仅使用大量的正常样本,而且可以使用少量的故障样本对决策边界进一步修订,其基本思想是通过训练构造一个包围正常样本的超球,在使超球体积最小化的同时,进一步使超球边界和故障样本之间的间隔最大化。本文首先使用仿真数据进行不平衡样本下的分类实验,直观地表明了小球大间隔方法在不平衡样

收稿日期:2011-07-25

基金项目:国家自然科学基金资助项目(61179057);航空科学基金资助项目(2007ZB52022)

本学习下的优越性;然后将该方法应用到滚动轴承故障检测中,并将其与传统的支持向量机和支持向量数据描述进行了比较,验证了该方法在解决不平衡样本下机械故障检测问题中的优越性。

### 1 小球大间隔方法

小球大间隔方法<sup>[14]</sup>是一种针对训练中拥有大量正常样本和少量异常样本情况的异常检测方法,其集成了一类分类方法(支持向量数据描述)和传统两类分类方法(支持向量机)的思想。一方面,与支持向量数据描述类似,小球大间隔方法通过在特征空间中构造一个包围正常样本的超球来进行异常检测,若一个测试样本落入超球内部,则将其分类为正常,否则,将其分类为异常。为了减小将异常样本分类为正常样本的可能性,该超球的体积被最小化。另一方面,受支持向量机大间隔思想的启发,为了进一步减小将异常样本分类为正常样本的可能性,小球大间隔方法要求超球边界与异常样本之间的间隔最大化。

给定  $m_1$  个正常样本  $(x_1, y_1), (x_2, y_2), \dots, (x_{m_1}, y_{m_1})$  和  $m_2$  个异常样本  $(x_{m_1+1}, y_{m_1+1}), (x_{m_1+2}, y_{m_1+2}), \dots, (x_{m_1+m_2}, y_{m_1+m_2})$ , 其中  $x_k \in \mathbf{R}^d$  为输入样本( $k = 1, 2, \dots, m_1 + m_2$ ),  $y_k$  为样本类别标号,当  $i = 1, 2, \dots, m_1$  时,  $y_i = 1$ , 当  $j = m_1 + 1, m_1 + 2, \dots, m_1 + m_2$  时,  $y_j = -1$ , 令  $n = m_1 + m_2$  ( $m_2$  远小于  $m_1$ ), 则小球大间隔方法可以表示为以下最优化问题:

$$\min_{R, c, \rho, \xi} R^2 - \nu \rho^2 + \frac{1}{\nu_1 m_1} \sum_{i=1}^{m_1} \xi_i + \frac{1}{\nu_2 m_2} \sum_{j=m_1+1}^n \xi_j \quad (1)$$

$$\text{s. t.} \quad \|\varphi(x_i) - c\|^2 \leq R^2 + \xi_i \quad (2)$$

$$\|\varphi(x_j) - c\|^2 \geq R^2 + \rho^2 - \xi_j \quad (3)$$

$$\xi_k \geq 0 \quad (4)$$

式中,  $\varphi(x_i), \varphi(x_j)$  分别为正常样本和异常样本在特征空间中的位置;  $c, R$  分别为在特征空间中建立的超球的球心位置和半径;  $\rho^2$  为超球边界与异常样本之间的间隔;  $\xi$  为松弛向量,  $\xi = (\xi_1, \xi_2, \dots, \xi_n) \in \mathbf{R}^n$ ;  $\nu, \nu_1, \nu_2$  为三个正常数。

根据上述最优化问题,最小化目标函数将使超球的半径  $R$  最小化,同时超球边界与异常样本之间的间隔  $\rho^2$  最大化,因此将该异常检测方法称为小球大间隔方法。

为了导出式(1)~式(4)的对偶形式,定义以下 Lagrange 函数:

$$L(R, c, \rho, \xi, \alpha, \beta) = R^2 - \nu \rho^2 +$$

$$\frac{1}{\nu_1 m_1} \sum_{i=1}^{m_1} \xi_i + \frac{1}{\nu_2 m_2} \sum_{j=m_1+1}^n \xi_j - \sum_{k=1}^n \beta_k \xi_k +$$

$$\sum_{i=1}^{m_1} \alpha_i (\|\varphi(x_i) - c\|^2 - R^2 - \xi_i) - \sum_{j=m_1+1}^n \alpha_j (\|\varphi(x_j) - c\|^2 - R^2 - \rho^2 + \xi_j) \quad (5)$$

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$$

$$\beta = (\beta_1, \beta_2, \dots, \beta_n)$$

其中,  $\alpha_i, \beta_i$  为 Lagrange 乘子,  $\alpha_i \geq 0, \beta_i \geq 0$ 。令  $L(R, c, \rho, \xi, \alpha, \beta)$  关于原始变量的导数为零,可得

$$\frac{\partial L}{\partial R} = 2R(1 - \sum_{i=1}^n \alpha_i y_i) = 0 \quad (6)$$

$$\frac{\partial L}{\partial \rho} = 2\rho(\sum_{j=m_1+1}^n \alpha_j - \nu) = 0 \quad (7)$$

$$\frac{\partial L}{\partial \xi_i} = \frac{1}{\nu_1 m_1} - \alpha_i - \beta_i = 0 \quad (8)$$

$$\frac{\partial L}{\partial \xi_j} = \frac{1}{\nu_2 m_2} - \alpha_j - \beta_j = 0 \quad (9)$$

$$\frac{\partial L}{\partial c} = 2c \sum_{i=1}^n \alpha_i y_i - 2 \sum_{i=1}^n \alpha_i y_i \varphi(x_i) = \mathbf{0} \quad (10)$$

由式(10)、式(6)可得

$$c = \sum_{i=1}^n \alpha_i y_i \varphi(x_i) / \sum_{i=1}^n \alpha_i y_i = \sum_{i=1}^n \alpha_i y_i \varphi(x_i) \quad (11)$$

将式(6)~式(9)和式(11)代入式(5),即可得到上文最优化问题的对偶形式:

$$\min_{\alpha \in \mathbf{R}^n} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \sum_{i=1}^n \alpha_i y_i K(x_i, x_i) \quad (12)$$

$$\text{s. t.} \quad 0 \leq \alpha_i \leq \frac{1}{\nu_1 m_1} \quad (13)$$

$$0 \leq \alpha_j \leq \frac{1}{\nu_2 m_2} \quad (14)$$

$$\sum_{i=1}^n \alpha_i y_i = 1 \quad (15)$$

$$\sum_{i=1}^n \alpha_i = 2\nu + 1 \quad (16)$$

求解出对偶问题(式(12)~式(16))后,为了计算半径  $R$ ,考虑下列集合:

$$S = \{x_i \mid 0 < \alpha_i < \frac{1}{\nu_1 m_1}\}$$

根据 KKT(Karush-Kuhn-Tucker) 条件,对于  $S$  中的样本,式(2)和式(3)中的等号成立,同时松弛变量为零。令  $n = |S|$ , 则

$$R^2 = P/n \quad (17)$$

$$P = \sum_{x_i \in S} \|\varphi(x_i) - c\|^2 =$$

$$\sum_{x_i \in S} (K(x_i, x_i) - 2 \sum_{k=1}^n \alpha_k y_k K(x_i, x_k) + \langle c, c \rangle) \quad (18)$$

根据式(11),可求得

$$\langle c, c \rangle = \langle \sum_{i=1}^n \alpha_i y_i \varphi(x_i), \sum_{j=1}^n \alpha_j y_j \varphi(x_j) \rangle = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad (19)$$

为了对一个测试样本  $x \in \mathbf{R}^d$  进行分类,只须根据决策函数判断该样本是否在训练阶段构造的

超球体内。决策函数的表达式为

$$f(x) = \text{sgn}(R^2 - \|\varphi(x) - c\|^2) = \text{sgn}(R^2 - \langle c, c \rangle - K(x, x) + 2 \sum_{k=1}^n \alpha_k y_k K(x, x_k)) \quad (20)$$

## 2 小球大间隔方法在不平衡样本学习下的优越性比较与分析

为了直观地比较支持向量机、支持向量数据描述以及小球大间隔方法应用于不平衡样本时的局限性与优越性,利用二维仿真数据对三种方法进行了不平衡样本下的分类实验。

仿真数据通过随机均匀分布产生,具体产生办法为:在由横坐标 $[0, 1]$ 和纵坐标 $[0, 1]$ 形成的区域内根据均匀分布随机产生 200 个正类训练样本,在由横坐标 $[1, 2]$ 和纵坐标 $[0, 1]$ 形成的区域内根据均匀分布随机产生 20 个负类训练样本,训练样本中正类样本和负类样本的不平衡比为 10 : 1,然后采用相同的方法另外产生 100 个正类

样本和 100 个负类样本用于测试,实验数据的具体细节如表 1 所示。

表 1 仿真实验数据

方法	不平衡比	训练样本个数		测试样本个数	
		正类	负类	正类	负类
SVM	10 : 1	200	20	100	100
SVDD		200		100	100
SSLM	10 : 1	200	20	100	100

3 种方法的核函数均选取高斯核函数。SVM 和 SSLM 的参数通过五折交叉验证进行选取。由于 SVDD 仅使用一类样本进行训练,因此不适合使用交叉验证的方法选取参数,鉴于该方法和 SSLM 一样都是通过构造一个包围正类样本的超球来进行分类,因此,为公平起见,选取和 SSLM 一样的核参数,另一个惩罚参数选取 1,即要求在训练集上没有误分。使用上述参数选取方法选取参数后进行分类实验,3 种方法在训练集和测试集上的分类结果分别如图 1 和图 2 所示,具体的识别率如表 2 所示。

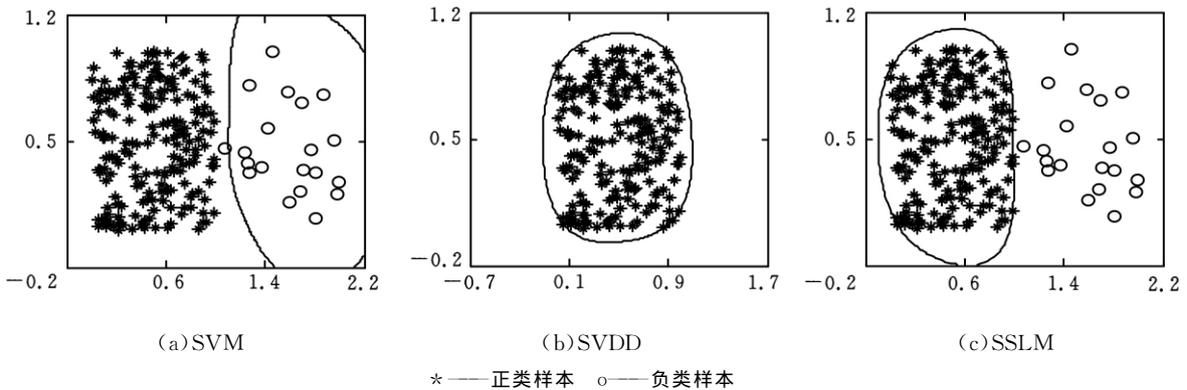


图 1 三种方法在训练集上的分类结果

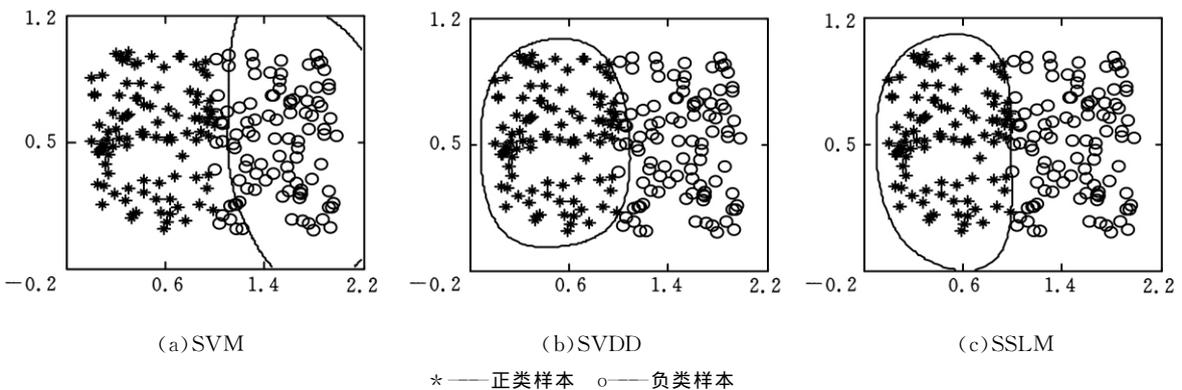


图 2 三种方法在测试集上的分类结果

建立在结构风险最小化之上的 SVM,通过在模型的复杂性和训练误差之间寻求折中,而不是一味地追求训练误差最小化,从而能够有效地避

免过拟合现象,表现出良好的推广能力。但是,当训练样本严重不平衡时,一方面较少一类的样本很容易远离理想的分类面,另一方面 SVM 软间

表 2 三种方法在仿真数据上的识别率

方法	不平衡比	训练集		测试集	
		正类 识别率	负类 识别率	正类 识别率	负类 识别率
SVM	10 : 1	1.00	0.95	1.00	0.81
SVDD		1.00		1.00	0.93
SSLM	10 : 1	1.00	1.00	0.98	1.00

隔的特点使得 SVM 训练得到的分类面会向样本较少一类偏移<sup>[15]</sup>,当使用该分类面对测试样本进行分类时,样本个数较少的一类会具有较高的误识率。从图 1a 可以看出:由于正类样本较多,因此在理想分类面附近分布有很多样本;由于负类样本较少,仅有一个样本接近理想分类面。这种情况下,SVM 在软间隔特点的作用下,为了获得更大的间隔,最终训练得到的分类面向负类样本方向发生了偏移,并越过了一个负类样本,导致在训练集上该类样本有一个发生了误分。总之,由于上述两个原因,SVM 在训练样本不平衡的情况下得到的分类面明显向样本较少的负类方向发生了偏移。从图 2a 可以看出,当使用该分类面对测试样本进行分类,正类样本没有误分,其识别率达到了 100%,而负类样本有相当一部分发生了误分,其识别率仅为 81%。

SVDD 作为一种一类分类方法,在训练中只须使用正常样本,其基本思想是通过构造一个包围正常样本的超球来进行异常检测,若测试样本落入超球内部,则将其分类为正常,否则将其分类为异常。通过引入核技巧,该方法可以获得灵活的描述边界。从该方法的原理可以看出,SVDD 进行异常检测的效果取决于得到的描述边界是否紧凑,若描述边界非常紧凑,则异常样本很难落入超球内部,从而可以获得较好的异常检测效果,否则,异常检测的效果会较差。但是,只有当核参数选取恰当时,SVDD 才能获得紧凑的描述边界,由于该方法在训练时仅使用了正常样本,当训练样本位于高维空间时,仅通过正常样本很难判断选取的核参数对应的描述边界是否紧凑,若训练得到的描述边界不是非常紧凑,使用该方法进行异常检测时接受异常样本的风险会较高。另外,由于 SVDD 要求包围正常样本的超球尽可能小,这使得其和 SVM 相比,正常样本的识别率容易偏低。图 1b 所示为 SVDD 使用正类样本训练得到的分类面,该分类面包围住了所有的正类样本,即其在训练集上的识别率达到了 100%,但是从图 2b 可以看出,使用该分类面对测试样本进行分类时,正类样本没有一个发生误分,负类样本有 7 个发生了误分,即正类样本的识别率为 100%,负类

样本的识别率为 93%。

SSLM 也是通过构造一个包围正常样本的超球来进行异常检测,这一点与 SVDD 类似,不同点在于该方法在训练中引入了异常样本,在最小化超球的同时,进一步使超球边界和异常样本之间的间隔最大化,因此,与 SVDD 相比,该方法一般可以获得更加紧凑的描述边界,从而能够降低接受异常样本的风险。由于该方法同样要求包围正常样本的超球尽可能小,因此和 SVM 相比,其正常样本的识别率也容易偏低。图 1c 所示为 SSLM 通过训练得到的分类面,其对正类样本和负类样本的识别率均达到了 100%。对比 SSLM 和 SVDD 的分类面,可以看出,在靠近负类样本一侧,SSLM 的分类面更加紧凑,这一点正是通过大间隔的要求获得的。图 2c 所示为 SSLM 在测试样本上的分类结果,其对正类样本的识别率为 98%,略低于另外两种方法的正类识别率,对负类样本的识别率为 100%,明显高于另外两种方法的负类识别率。

根据以上实验和分析可以看出,当训练样本高度不平衡时,SVM 的分类面会偏向样本较少的类,这使得其对样本较少类的识别率容易偏低,因此训练样本不平衡会明显降低 SVM 的性能;SVDD 仅使用一类样本进行训练,虽不存在训练样本不平衡的问题,但该特点导致其核参数选择比较困难,不能保证对应的描述边界一定非常紧凑,若核参数选择不当,该方法对异常样本的识别率容易偏低;SSLM 在训练中同时使用正常样本和异常样本,通过构造一个包围正常样本的超球来进行异常检测,这一特点使得其面对不平衡的训练样本时不存在 SVM 的缺点,而大间隔的特点又克服了 SVDD 的不足,因此 SSLM 可以作为一种很好的不平衡样本下的异常检测方法。机械故障检测是一种典型的不平衡样本下的异常检测问题,因此 SSLM 可以用于解决不平衡样本下的机械故障检测问题。

### 3 滚动轴承故障检测应用

为验证 SSLM 在不平衡样本下机械故障检测中的优越性,本文使用滚动轴承故障模拟实验台数据进行了不平衡样本下的故障检测实验,作为对比,同时使用 SVM 和 SVDD 进行了实验。滚动轴承故障模拟实验台如图 3 所示,其中,通道 1 的传感器用于测试滚动轴承转速,通道 2 和通道 4 的传感器用于测试水平加速度信号,通道 3 的传感器用于测试垂直加速度信号。采用 4 个

6304 型滚动轴承进行实验,其中,1 个为正常轴承,另外 3 个被设置有内圈故障、外圈故障和滚动体故障(故障通过电火花技术加工而成)。实验中轴承的转速为 1500r/min、1800r/min 和 2000 r/min,采样频率为 10kHz。根据上述实验条件采集 4 种状态(正常、内圈故障、外圈故障和滚动体故障)下的滚动轴承振动信号各 300 组数据,每一组数据包含 4096 个数据点。实验中将内圈故障、外圈故障和滚动体故障样本混合组成故障类样本。根据训练样本中正常样本和故障样本的不同比例(1 : 1、10 : 1、20 : 1、40 : 1),将实验分为 4 组,实验数据的具体细节如表 3 所示。

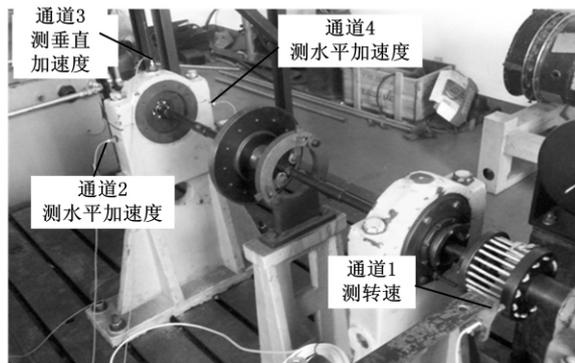


图 3 滚动轴承故障模拟实验台

表 3 滚动轴承故障检测实验数据

方法	不平衡比	训练样本个数		测试样本个数	
		正常	故障	正常	故障
SVM	1 : 1	200	200	100	100
	10 : 1	200	20	100	100
	20 : 1	200	10	100	100
	40 : 1	200	5	100	100
SVDD		200		100	100
SSLML	1 : 1	200	200	100	100
	10 : 1	200	20	100	100
	20 : 1	200	10	100	100
	40 : 1	200	5	100	100

本文将滚动轴承振动信号的波形指标、峰值指标、脉冲指标、裕度指标和峭度指标作为滚动轴承的故障特征,将这 5 个特征组成五维特征向量并作为分类器的输入。

传统的分类器性能评估注重总体性能,一般将分类正确的样本个数和测试样本总个数的比值作为评价指标。但是,当数据高度不平衡,并且不同类别的误分代价不一样时,这种方法存在很大的局限性,原因是假如有 100 个样本,其中,正类样本 99 个,负类样本 1 个,如果分类器将所有样本分为正类,根据上述方法,分类器的识别率为 99%,但是,由于负类的识别率为 0,如果故障对应负类,那么故障的识别率就为 0,对于故障检测来说,这样的结果没有任何意义。针对该问题,当

样本高度不平衡时,分类器的性能评估一般采用正类识别率和负类识别率的几何平均值<sup>[16]</sup>,即令  $a^+$ 、 $a^-$  分别表示正类识别率和负类识别率,则总识别率为  $\sqrt{a^+ a^-}$ 。在下面的故障检测实验中,将采用这种方法来评估 3 种方法的故障检测性能。

实验中,3 种方法的核函数统一选择高斯核函数。参数选取方法同上节仿真实验一样,即 SVM 和 SSLML 的参数通过五折交叉验证进行选取。为公平起见,SVDD 选取和 SSLML 一样的核参数,另一个惩罚参数选取 1,即要求在训练集上没有误分。每一组实验中,训练样本和测试样本随机划分 10 次,对每一次划分的数据进行分类实验,然后将 10 次分类结果的平均值作为最终结果。实验结果如表 4 所示。

表 4 3 种方法在滚动轴承实验数据上的识别率

方法	不平衡比	训练集			测试集		
		正常识别率	故障识别率	总识 别率	正常识别率	故障识别率	总识 别率
SVM	1 : 1	1.0000	0.9945	0.9972	0.9840	0.9660	0.9750
	10 : 1	1.0000	0.9800	0.9899	1.0000	0.8670	0.9277
	20 : 1	1.0000	0.9100	0.9539	1.0000	0.6680	0.8173
	40 : 1	1.0000	0.8800	0.9381	1.0000	0.5250	0.7246
SVDD		1.0000		1.0000	0.9550	0.9110	0.9327
SSLML	1 : 1	1.0000	1.0000	1.0000	0.9410	1.0000	0.9701
	10 : 1	0.9990	1.0000	0.9995	0.9530	0.9810	0.9669
	20 : 1	0.9995	1.0000	0.9997	0.9550	0.9720	0.9635
	40 : 1	1.0000	1.0000	1.0000	0.9580	0.9660	0.9620

从表 4 中 SVM 的检测结果可以看出,当训练样本平衡时,在训练集和测试集上,SVM 对于正常样本和故障样本均取得了较高的识别率,这表明当训练样本平衡时,SVM 是一种很好的故障检测方法。当训练样本不平衡时,在训练集上,SVM 对正常样本的识别率达到了 100%,对故障样本的识别率较训练样本平衡时出现了一定下降,而且不平衡比越高,对故障样本的识别率下降越明显;在测试集上,与训练样本平衡时相比,SVM 对正常样本的识别率稍有上升,达到了 100%,但对故障样本的识别率出现了明显下降,而且随着不平衡比的增加,对故障样本的识别率急剧下降,这表明在训练样本不平衡时,SVM 对正常样本的识别率容易偏高,而对故障样本的识别率容易偏低。在机械故障检测中,由于将故障误判为正常的代价远高于将正常误判为故障的代价,因此一般希望对故障的识别率能够较高,但是 SVM 在训练样本不平衡情况下的性能正好与该目标相反。

从表 4 中 SVDD 的检测结果可以看出,由于该方法仅使用正常样本进行训练,因此不存在在训

训练样本不平衡的问题。在训练集上,其对正常样本的识别率达到了100%;在测试集上,其对正常样本的识别率低于SVM对正常样本的识别率,对故障样本的识别率低于训练样本平衡时SVM对故障样本的识别率,但高于训练样本不平衡时SVM对故障样本的识别率。总体来说,当训练样本不平衡时,SVDD的故障检测性能优于SVM的故障检测性能。

从表4中SSLM的检测结果可以看出,在训练集上,对于各种不平衡比,其对正常样本和故障样本的识别率都接近或达到了100%;在测试集上,该方法对正常样本的识别率略低于SVM对正常样本的识别率,与SVDD对正常样本的识别率基本相当,但是对故障样本的识别率明显高于SVM和SVDD对故障样本的识别率。总体来说,对于各种不平衡比,在训练集和测试集上,SSLM对于正常样本和故障样本的识别率均取得了较大的值,这表明SSLM基本不受训练样本不平衡的影响。此外还可以看出,SSLM对故障样本的识别率较正常样本更高一些,这一点与故障检测中更加重视故障识别率的目标一致。总之,SSLM基本不受训练样本不平衡的影响和其更加注重故障识别率的特点使得其非常适合不平衡样本下的机械故障检测。

#### 4 结论

(1)当训练样本严重不平衡时,支持向量机训练得到的分类面会向训练样本较少的故障类方向偏移,从而会导致故障类具有较高的误识率,而且训练样本不平衡程度越严重,支持向量机对故障类的误识率越高。

(2)支持向量数据描述在训练中仅使用了正常样本,该特点导致其核参数选取困难,不能保证获得的描述边界一定非常紧凑,若核参数选取不当,容易造成故障识别率偏低。

(3)小球大间隔方法在最小化超球的同时,进一步使超球边界和故障样本之间的间隔最大化,这使得其对故障的识别率能有很好的保证,可以作为解决不平衡样本下机械故障检测问题的有效方法。

#### 参考文献:

- [1] Vapnik V N. The Nature of Statistical Learning Theory[M]. New York:Springer-Verlag,1995.
- [2] 陈果. 基于遗传算法的支持向量机分类器模型参数优化[J]. 机械科学与技术,2007,26(3):347-350.

- [3] 尉询楷,陆波,汪诚,等. 支持向量机在航空发动机故障诊断中的应用[J]. 航空动力学报,2004,19(6):844-848.
- [4] 徐启华,师军. 应用SVM的发动机故障诊断若干问题研究[J]. 航空学报,2005,26(6):686-690.
- [5] 吴峰崎,孟光. 基于支持向量机的转子振动信号故障分类研究[J]. 振动工程学报,2006,19(2):238-241.
- [6] Yuan S F,Chu F L. Support Vector Machines-based Fault Diagnosis for Turbo-pump Rotor[J]. Mechanical Systems and Signal Processing,2006,20(4):939-952.
- [7] 唐浩,屈梁生. 基于支持向量机的发动机故障诊断[J]. 西安交通大学学报,2007,41(9):1121-1126.
- [8] Widodo A,Yang B S,Han T. Combination of Independent Component Analysis and Support Vector Machines for Intelligent Faults Diagnosis of Induction Motors[J]. Expert Systems with Applications,2007,32(2):299-312.
- [9] Tax D,Duin R. Support Vector Domain Description[J]. Pattern Recognition Letters,1999,20(11/13):1191-1199.
- [10] Tax D,Duin R. Support Vector Data Description[J]. Machine Learning,2004,54(1):45-66.
- [11] 李凌均,张周锁,何正嘉. 基于支持向量数据描述的机械故障诊断研究[J]. 西安交通大学学报,2003,37(9):910-913.
- [12] 王自营,邱绵浩,安钢,等. 基于一类超球面支持向量机的机械故障诊断研究[J]. 振动工程学报,2008,21(6):553-558.
- [13] 李强,王太勇,王正英,等. 基于EMD和支持向量数据描述的故障智能诊断[J]. 中国机械工程,2008,19(22):2718-2721.
- [14] Wu M R,Ye J P. A Small Sphere and Large Margin Approach for Novelty Detection Using Training Data with Outliers[J]. IEEE Transactions Pattern Analysis and Machine Intelligence,2009,31(11):2088-2092.
- [15] Akbani R,Kwek S,Japkowicz N. Applying Support Vector Machines to Imbalanced Datasets[C]//Proceedings of the 15th European Conference on Machine Learning. Pisa,Italy,2004:39-50.
- [16] He H B,Garcia E A. Learning from Imbalanced Data[J]. IEEE Transactions on Knowledge and Data Engineering,2009,21(9):1263-1284.

(编辑 张洋)

作者简介:郝腾飞,男,1983年生。南京航空航天大学民航学院博士研究生。研究方向为航空发动机状态监测与故障诊断。  
陈果,男,1972年生。南京航空航天大学民航学院教授、博士研究生导师。