

# 基于样本熵与决策树调节算法的轴承故障识别

李六杏, 周黄丽

(安徽经济管理学院信息工程系, 合肥 230031)

**摘要:** 轴承故障是导致机器发生事故的重要原因之一。为更好地识别出故障类型, 使用一种包络样本熵和决策树门限值自适应调节算法相结合的方法。首先将信号分解成若干 IMF 之和, 选取包含丰富故障信息的 IMF 求其包络信号的样本熵, 最后通过决策树自适应调节门限值准确判断出轴承故障类型。分析结果表明, 该方法不仅可以通过反馈减少运算量, 而且能够通过决策树门限值的自适应调节来提高轴承故障的识别率, 综合识别率可达到 96.75%。

**关键词:** 包络样本熵; 决策树; 轴承故障; 自适应调节; 门限值

中图分类号: TP206

文献标识码: A

文章编号: 1672-352X (2017)05-0936-05

## Based on sample entropy and decision tree algorithm for regulation of bearing fault diagnosis

LI Luxing, ZHOU Huangli

(Department of Information Engineering, Anhui Institute of Economic Management, Hefei 230031)

**Abstract:** Bearing failure is one of the important reasons for machine accidents. To better identify fault types, this paper used a method which combines envelope sample entropy with decision tree gate limit value adaptive adjustment algorithm. We firstly decomposed the signal into several IMFs and selected the sample entropy of the envelope signal from the IMF which contains a number of diagnostic information, and then, the type of bearing failure was accurately determined by decision tree adaptive adjustment gate limit value. Analysis results showed that this method can not only reduce the computational complexity through the feedback, but also improve the recognition rate of bearing failure through the decision tree threshold adaptive adjustment. Comprehensive recognition rate can reach 96.75%.

**Key words:** envelope sample entropy; decision tree; bearing failure; adaptive adjustment; gate limit value

轴承发生故障时产生的信号通常为非平稳信号, 如何将轴承故障的特征信息从非平稳信号中提取出来是识别轴承故障的关键所在。康海英等<sup>[1]</sup>提出了一种基于阶次跟踪和 HHT 瞬时相位法的轴承故障诊断方法, 时域里的非平稳信号转换为角域里的平稳信号。董红生等<sup>[2]</sup>提出一种基于 HHT 边际谱熵和能量谱熵的心率变异信号的分析方法, 为临床 HRV 信号及其他复杂生理信号的分析提供一种有效的分析方法。傅勤毅等<sup>[3]</sup>采用一种无频带错位的小波包算法进行滚动轴承的故障特征信号提取, 清晰地刻画出轴承故障冲击的特征函数。杨宇等<sup>[4]</sup>提出了一种基于经验模态分解和神经网络的滚动轴承故障诊断方法来识别轴承故障。但这些方法不仅算法复杂度高, 而且计算量较大。

本研究使用包络样本熵和决策树门限值自适应调节算法相结合的方法, 将 EMD 分解后得到的 IMF 分量经过提取后计算其包络样本熵, 使用包络样本熵训练决策树, 最后使用生成的决策树并自动调节门限值来提高轴承故障的识别率。该方法不仅算法复杂度小, 具有可靠性实用性, 而且大大降低了计算量。

## 1 材料与方法

### 1.1 材料

滚动轴承故障数据来源于西安交通大学机械工程学院智能仪器与监测诊断研究所, 本次研究采用的轴承包括状态良好的轴承、外圈有剥落的轴承、内圈有剥落的轴承以及滚珠有剥落的轴承。轴承型

收稿日期: 2017-03-19

基金项目: 安徽省高校自然科学研究重点项目 (KJ2015A394) 和安徽经济管理学院院级课题 (YJKT1516YB07) 共同资助。

作者简介: 李六杏, 副教授。E-mail: llxin001080@126.com

号及相关参数如下: 轴承型号 6308, 滚珠数  $Z=8$ , 钢球直径  $d=15$  mm, 滚道节径  $E=65.5$  mm, 接触角  $\alpha=6$  度。采样频率为 20 kHz, 每种故障的轴承信号采样 4 次。

## 1.2 方法

**1.2.1 包络样本熵** (1) EMD 理论。EMD 方法可以把原始信号分解成一系列本征模态函数 (IMF)。每个 IMF 都符合以下 2 个条件:

极值点与过零点数目相差不超一个; 上下包络线的局部平均值为 0。原始信号可以看成若干个满足上述条件的 IMF 之和, 每个 IMF 按照以下步骤确定<sup>[5-8]</sup>:

1) 先找到信号  $x(t)$  的全部极值点, 然后采用三次样条曲线分别对全部局部极大值点和局部极小值点进行拟合。拟合的结果分别称为上包络线和下包络线。上、下包络线包含了整个信号序列。

2) 令  $u_1$  为上、下包络线的平均值, 记

$$y_1(t) = x(t) - u_1 \quad (1)$$

3) 验证  $y_1(t)$  是不是 IMF。如果  $y_1(t)$  不符合 IMF 条件, 就把  $y_1(t)$  当作初始信号, 再次进行 1)、2) 步骤, 直到  $y_1(t)$  符合 IMF 条件。令  $y_1(t) = c_1(t)$ , 则得到  $x(t)$  的第一个同时也是频率最高的一个 IMF 分量。

4) 把  $c_1(t)$  从  $x(t)$  中剔除, 得到去除第一个 IMF 的剩余部分  $r_1(t)$ 。则

$$r_1(t) = x(t) - c_1(t) \quad (2)$$

把  $r_1(t)$  当作初始信号, 继续 1、2、3 步骤得到第二个 IMF 分量。重复进行  $n$  次, 获得  $n$  个 IMF 分量。当  $c_n(t)$  或  $r_n(t)$  是单调信号或者小于一定的值时, 分解结束。此时

$$x(t) = \sum_{i=1}^n c_i(t) + r_n(t) \quad (3)$$

其中,  $r_n(t)$  表示信号趋势的残余分量。经过以上步骤, 一个完整的信号就被分解为不同频率特征的 IMF 分量和残余分量之和的形式。

(2) IMF 包络信号。通过希尔伯特变换来获得分解后 IMF 分量的包络信号。一个实信号  $x(t)$  的希尔伯特变换定义为:

$$H[x(t)] = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{x(\tau)}{t-\tau} d\tau \quad (4)$$

具体地, 对 EMD 分解后的 IMF 分量其形式为

$$H[c_i(t)] = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{c_i(\tau)}{t-\tau} d\tau \quad (5)$$

对应的包络信号为

$$A(t) = \sqrt{c_i^2(t) + H^2[c_i(t)]} \quad (6)$$

(3) 包络样本熵。对于一个  $N$  点时间序列  $x(1), x(2), \dots, x(N)$ , 其样本熵计算步骤如下:

1) 由原时间序列构成一个  $m$  维向量

$$X_m(i) = \{x(i), x(i+1), \dots, x(i+m-1)\}, 1 \leq i \leq N-m+1 \quad (7)$$

2) 定义  $X_m(i)$  与  $X_m(j)$  之间的距离

$$d[X_m(i), X_m(j)] = \max(|x(i+k) - x(j+k)|), 1 \leq k \leq m-1 \quad (8)$$

3) 给定容限  $r$ , 统计每个  $X_m(i)$  对应的  $d[X_m(i), X_m(j)] \leq r$  的数目, 记为  $A_i$ , 把  $A_i$  与  $N-M+1$  的比值记为

$$B_i^m(r) = \frac{A_i}{N-m+1}, 1 \leq i \leq N-m \quad (9)$$

4) 求出  $B_i^m(r)$  的平均值  $B^m(r)$ ,

$$B^m(r) = \frac{1}{N-m} \sum_{i=1}^{N-m} B_i^m(r) \quad (10)$$

5) 用同样的方法求出  $B_i^{m+1}(r)$ , 该时间序列理论上的样本熵定义为

$$SanpEn(m, r) = \sum_{N \rightarrow \infty} -\ln \frac{B^{m+1}(r)}{B^m(r)} \quad (11)$$

当  $N$  为有限值时, 样本熵的估计值为

$$SanpEn(m, r, N) = -\ln \frac{B^{m+1}(r)}{B^m(r)} \quad (12)$$

由上面对样本熵的计算可以看出, 参数  $m$  和参数  $r$  对样本熵的计算结果影响很大。根据相关资料<sup>[9]</sup>,  $m=1$  或者  $2$ ,  $r=0.1 \sim 0.25SD$  ( $SD$  为原始时间序列的标准差) 时计算结果比较好。选用  $m=2$ ,  $r=0.15SD$ 。

**1.2.2 决策树** 决策树<sup>[10]</sup>是应用广泛的分类算法之一, 核心算法是 C4.5。C4.5 算法由 Quinlan 在 ID3 的基础上提出的<sup>[11]</sup>。ID3 算法用来构造决策树。决策树是一种类似流程图的树结构, 其中每个内部节点 (非树叶节点) 表示在一个属性上的测试, 每个分枝代表一个测试输出, 而每个树叶节点存放一个类标号。一旦建立好了决策树, 对于一个未给定类标号的元组, 跟踪一条有根节点到叶节点的路径, 该叶节点就存放着该元组的预测。决策树的优势在于不需要任何领域知识或参数设置, 适合于探测性的知识发现。决策树是以实例为基础的归纳学习算法, 对一组无次序无规则的样本推理出决策树表示分类规则。对目标类尝试进行最佳的分割, 从跟到叶子节点都有一条路径, 也即是一条规则。决策树的优点: 计算量不大; 可以处理连续和种类字段; 清晰的表现出重要特征属性, 并对噪声有很好的健壮性。

决策树分为树的构造修剪两部分算法。构造的

关键部分就是选择好的逻辑判断和属性。C4.5 是用信息增益率来选择属性，克服了用信息增益来选择属性时偏向选择值多的属性的不足。信息熵定义为：

$$SplitInmation(S,A) = -\sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|} \quad (13)$$

其中， $S$  是样本集， $S_i$  是  $S$  的子集，按照属性  $a$  分类的信息增益为  $Gain(S,A)$ ，则信息增益率为  $GainRatio(S,A)$ 。C4.5 决策树的每个节点是使用信息增益比来度量属性的选择，选择信息增益比最大的属性作为当前节点属性。

$$Gain(S,A) = SplitInmation(S,A) - SplitInmation(S_a,A) \quad (14)$$

$$GainRatio(S,A) = \frac{Gain(S,A)}{SplitInmation(S,A)} \quad (15)$$

**1.2.3 诊断算法流程** 包络样本熵和决策树门限值自适应调节算法主要包括信号处理环节、决策树训练以及信息反馈这 3 个步骤：

(1) 对原始数据进行 EMD 分解，得到若干个 IMF 分量，并且对各个 IMF 分量做自相关分析，选取其中和原始数据相关性较大的前几个 IMF 分量。

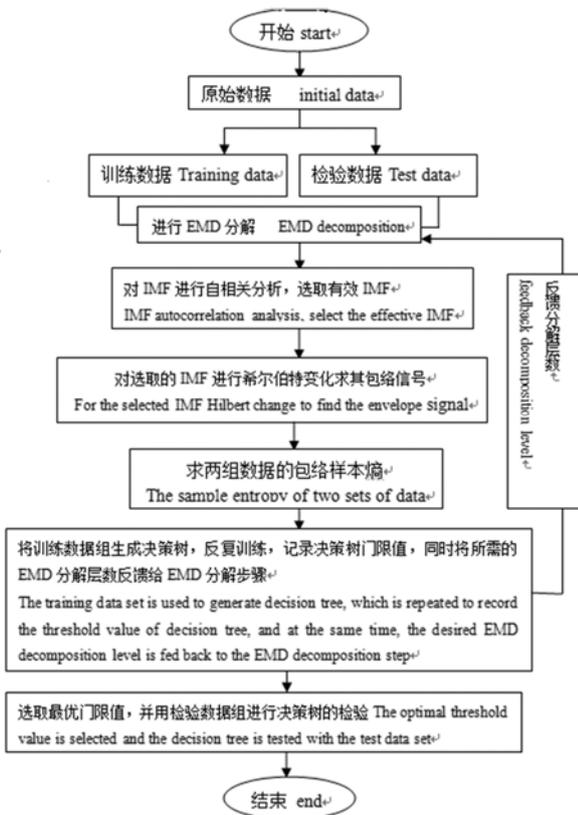


图 1 算法流程

Figure 1 Algorithm flowchart

(2) 对选取的 IMF 分量进行希尔伯特变换求其包络信号，然后对包络信号求样本熵。

(3) 训练决策树。由训练组数据生成决策树，调节训练组中各种故障所占的比例，反复训练决策树，将得到的门限值保存在相应的数组中，算法自动从中选取最优门限值，以使得识别率达到最高。

(4) 信息反馈。将决策树所需要的分解层数反馈给 EMD 分解环节，使之分解到一定的层数即可停止分解，减少运算量。

算法的流程图如图 1 所示。

## 2 结果与分析

### 2.1 决策树训练结果

将原始数据以 1 200 点长度为一个样本，每种故障取 70 个样本作为训练数据组，取 100 个样本作为检验数据组。分别计算两组数据的包络样本熵，用训练数据组生成并训练决策树，用检验数据组检验决策树。原始信号的时域波形图见图 2。

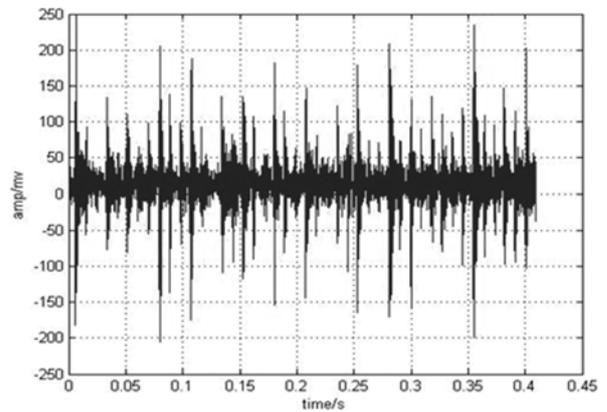


图 2 滚动体时域波形

Figure 2 Rolling time domain waveform

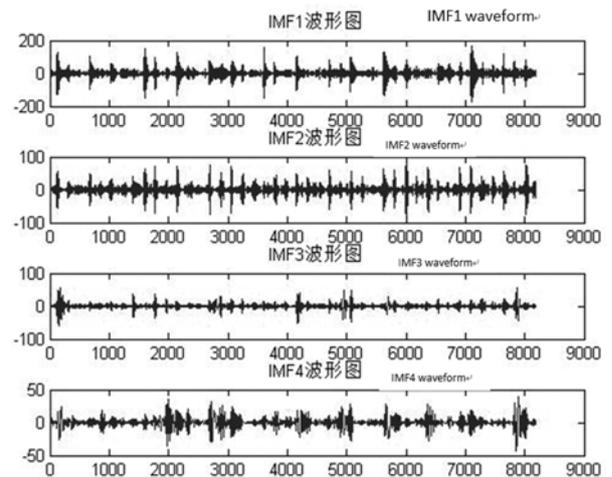


图 3 滚动体故障的前 4 个 IMF 波形

Figure 3 Rollover failure of the first four IMF waveforms

对原始信号进行 EMD 分解后提取有效的 IMF 分量。提取前 4 个 IMF 分量。将选取的 IMF 分量

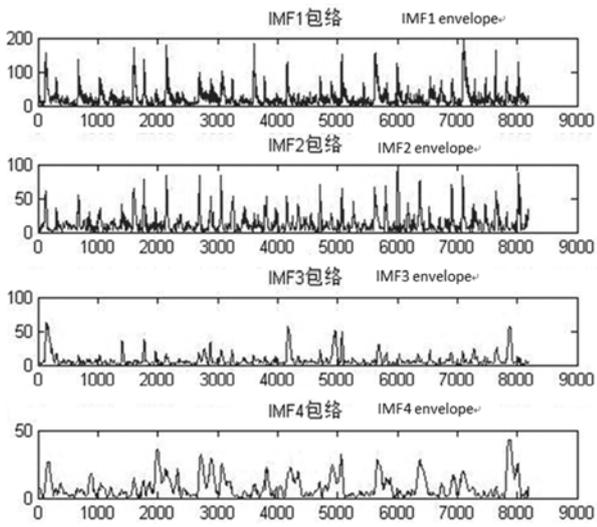


图 4 滚动体故障前 4 个 IMF 包络  
Figure 4 Rollover failure the first four IMF envelopes

做希尔伯特变换, 即可得到这些 IMF 分量的包络信号。选取的 IMF 分量的波形图及包络图如图 3、图 4 所示。

求出各中故障的包络样本熵后, 使用训练数据组中的样本熵生成 C4.5 决策树, 并通过自适应调节算法选择最优门限值, 决策树训练生成结果见图 5。图 5 中并没有 IMF3 和 IMF4 分量的包络样本熵, 说明只需要计算出前 2 个 IMF 分量的包络样本熵即可识别出所有的轴承故障类型, 将这些信息反馈给 EMD 分解过程, 使之分解到第 2 层时停止分解, 可以大大降低计算机的运算量。

分别求出外圈故障、滚动体故障轴承、内圈故障轴承和正常轴承数据的前 4 个 IMF 分量的包络信号样本熵, 样本熵值如下表 1, 可见由于内圈故障和外圈故障的 IMF 包络样本熵值较为接近, 因此该算法对内圈故障和外圈故障的识别率较低。

表 1 各种故障的前 4 个 IMF 分量包络样本熵

Table 1 The first four IMF components of various failures Envelope sample entropy

故障类型 Fault type	IMF1	IMF2	IMF3	IMF4
外圈故障 Outer fault	0.2346	0.0598	0.0353	0.0841
滚动体故障 Rolling fault	0.5024	0.2790	0.1243	0.1049
内圈故障 Inner fault	0.2268	0.2030	0.1599	0.1126
正常 Normal	1.1960	0.5263	0.4029	0.1577

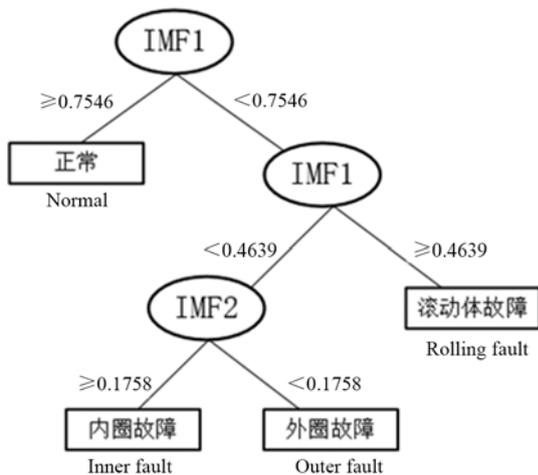


图 5 决策树训练结果  
Figure 5 Decision tree training results

### 2.2 检验结果与分析

使用检验数据组的信号包络样本熵对生成的决策树进行检验, 每种故障有 100 个样本, 检验结果如表 2 所示, 可见使用包络样本熵和决策树门限值自适应调节算法能够有效识别出轴承的故障类型。其中, 正常轴承的识别率达到 99%, 滚动体故障的识别率达到 98%。整体上的综合识别率可达 96.75%。

由表 3 可以看出, 采用了决策树门限值自适应调节算法的轴承故障识别率比没有使用门限自适应调节算法的识别率高出 1.75%, 比 SVM 和神经网络等方法的识别率都要高, 因此使用决策树门限值自适应调节算法能够较好的提高故障的识别率。

表 2 决策树检验结果

Table 2 Decision tree test results

故障类型 Fault type	外圈故障 Outer fault	滚动体故障 Rolling fault	内圈故障 Inner fault	正常 Normal
训练样本 Training sample	70	70	70	70
检验样本 Test sample	100	100	100	100
正确诊断 Accurate diagnosis	96	98	94	99
准确率 Accuracy rate	96%	98%	94%	99%

表 3 各种方法识别率的比较

Table 3 Comparison of the recognition rates of the various methods

方法名称 Method name	综合识别率 Outer fault %
决策树自适应调节算法 Adaptive adjustment algorithm of decision tree	96.75
普通的决策树算法 Common decision tree algorithm	95.00
SVM	94.00
神经网络 Neural network	92.00

### 3 结论

为了能够有效的识别出轴承故障类型, 提出将 IMF 包络样本熵和决策树门限值自适应调节算法相结合的方法。该方法能够通过决策树门限值的自适应调节来提高轴承故障的识别率, 并且可以通过反馈环节减少 EMD 的分解层数, 大大降低了运算量。相对于 SVM 和神经网络来说, 具有识别率高、运算量小、复杂度低的特点。

分析结果表明, 该方法能够有效识别出轴承故障类型, 相对于没有使用门限值自适应调节算法的决策树以及 SVM、神经网络等方法识别率有明显提高, 整体识别率可达到 96.75%。

### 参考文献:

- [1] 康海英, 栾军英, 郑海起, 等. 基于阶次跟踪和 HHT 瞬时相位法的轴承故障诊断[J]. 数据采集与处理, 2007, 22(1): 110-114.
- [2] 董红生, 邱天爽, 张爱华, 等. 基于 HHT 边际谱熵和能量谱熵的心率变异信号的分析方法[J]. 中国生物医学工程学报, 2010, 29(3): 336-344.
- [3] 傅勤毅, 章易程, 应力军, 等. 滚动轴承故障特征的小波提取方法[J]. 机械工程学报, 2001, 37(2): 30-32.
- [4] 杨宇, 于德介, 程军圣. 基于 EMD 与神经网络的滚动轴承故障诊断方法[J]. 振动与冲击, 2005, 24(1): 85-88.
- [5] HUANG N E, SHEN Z, LONG S R. A new view of nonlinear water waves: the Hilbert spectrum[J]. Annu Rev Fluid Mech, 1999, 31(1): 417-457.
- [6] HUANG N E, SHEN Z, LONG S R, et al. The empirical mode decomposition and the Hilbert spectrum for nonlinear and non-stationary time series analysis[C]// Proceedings of the Royal Society of London A: mathematical, physical and engineering sciences. London: The Royal Society, 1998, 454(1971): 903-995.
- [7] 胡荣华, 楼佩煌, 唐敦兵, 等. 基于 EMD 和免疫参数自适应 SVM 的滚动轴承故障诊断[J]. 计算机集成制造系统, 2013, 19(2): 438-447.
- [8] 王灿, 王嘉乐, 会强, 等. Hilbert-Huang 变换在机车滚动轴承故障诊断中的应用[J]. 振动, 测试与诊断, 2013, 33(增刊 1): 184-188.
- [9] PINCUS S M. Assessing serial irregularity and its implications for health[J]. ANN NY ACAD SCI, 2001, 954(1): 245-267.
- [10] JIA G F J, YUAN S F Y, TANG C W, et al. Fault diagnosis of roller bearing using feedback EMD and decision tree[C]//Electric Information and Control Engineering (ICEICE), Wuhan: 2011 International Conference on IEEE, 2011: 4212-4215.
- [11] Quinlan J R. C4.5: programs for machine learning[M]. Burlington: Morgan Kaufmann Publishers, 1993.