

一种新型复杂时间序列实时预测模型研究

王 军, 彭喜元, 彭 宇

(哈尔滨工业大学自动化测试与控制系, 黑龙江哈尔滨 150080)

摘 要: 针对复杂时间序列难以使用单一预测方法进行有效预测的问题, 本文提出一种新型多分辨率增量预测模型. 该模型首先使用经验模式分解方法对复杂时间序列分解, 然后对各分量分别进行增量核空间独立向量组合预测建模, 最后对各个分量预测结果等权求和集成为综合预测结果. 该预测模型可以实现对复杂时间序列的快速实时预测, 实验结果显示该模型在复杂时间序列预测上有良好的性能.

关键词: 复杂时间序列; 预测; 经验模式分解

中图分类号: TP391

文献标识码: A

文章编号: 0372-2112 (2006) 12A-2391-04

A Novel Real Time Predictor for Complex Time Series

WANG Jun, PENG Xi yuan, PENG Yu

(Department of Automatic Test and Control, Harbin Institute of Technology, Harbin, Heilongjiang 150080, China)

Abstract: The task of complex time series predicting is hard to be accomplished with only one single predicting model. In this paper, a novel multi-scale incremental predictor is proposed. This predictor decomposes the complex time series into a series of intrinsic mode functions (IMF) and a residual signal with empirical mode decomposition firstly, and then an Incremental Independent Vector Combination Predicting algorithm in Kernel Space (IIVCPKS) is constructed for predicting every IMF or residual signal. The proposed predictor is competent for predicting the complex time series in real time. Experimental results showed that the proposed method performed very well in the task of predicting complex time series.

Key words: complex time series; predicting; empirical mode decomposition (EMD)

1 引言

时间序列数据在一些新的数据库应用以及数据挖掘等领域中日趋重要. 而且这些数据往往是非平稳、非线性随机时间序列, 如气象数据、太阳黑子数据、激光数据、股票价格数据、网络流量数据、电力需求数据等. 时间序列的复杂行为可以用内在非线性动力系统描述, 但如何识别这样的非线性动力系统是个难题. 目前已经有很多进行时间序列预测的方法, 其中最为常用的预测方法为线性模型方法, 如 AR、ARMA、RLS 等, 其优点是模型简单、容易识别, 但应用线性模型预测非线性时间序列很难取得良好的预测效果^[1, 2]. 非线性模型由于其本身非线性特性适用于非线性时间序列建模, 如神经网络、模糊逻辑模型、双线性自回归模型、支持向量回归模型等^[3, 4]. 现有的预测方法大多采用单一的预测模型进行复杂时间序列建模, 但一个变化异常复杂的非线性、非平稳随机信号难以使用单一的预测模型进行^[5]. 此外, 复杂时间序列潜在物理规律是不断变化的, 所构造的预测模型应该跟踪这种变化, 这需要不断地对预测模型进行学习, 而现有的大部分非线性预测模型训练的时间复杂度很高, 所以非线性预测模型在线算法成为

近年来的一个研究热点.

本文将多分辨率经验模式分解与非线性预测模型的在线算法相结合提出一种新型预测模型, 该模型首先使用经验模式分解方法对复杂时间序列进行 n 阶分解, 然后分别对 $n-1$ 个本征模式分量和一个余项分量进行增量核空间独立向量组合预测建模, 并利用各个分量参数之间的关系减轻正交验证求参数的次数, 最后将这些分量预测结果等权求和集成为综合预测结果. 通过对 4 个真实复杂时间序列进行预测的结果显示, 本方法具有良好的预测性能.

2 经验模式分解

因为自然物理过程大多数是非线性和非平稳的, 这就限制了处理这些数据方法的可选择性. 现有的数据处理方法要么是针对线性非平稳过程的处理方法, 如小波变换、Wagner-Ville 分布和短时 Fourier 变换, 要么是针对非线性平稳且统计确定过程的处理方法, 如相空间表示法和时间延迟嵌入方法等. 小波和短时 Fourier 变换等方法是基于 Fourier 谱分析的, 通过可调时频窗函数来进行非平稳信号的时频分解, 在处理非平稳信号时可以达到一定的效果. 但使用谐波来表示非

线性信号时会发生能量向高频泄露,产生虚假频谱成分.相空间表示法和时间延迟嵌入方法在高噪声情况下将无法重构与真实物理过程微分同胚的动力系统.

由于上述这些方法很难处理现实世界中非线性、非平稳随机信号,人们迫切需求一种新型处理方法.1998年,Norden E.Huang提出一种适用于分析和处理非线性、非平稳随机信号的新方法——HHHT变换(Hilbert Huang Transform, HHT)^[7]. HHT变换由两个步组成:第一步将任意信号分解为若干本征模式分量(Intrinsic Mode Function, IMF)和一个余项,该步骤称为经验模式分解法(Empirical Mode Decomposition, EMD);第二步对每个本征模式或余项进行希尔伯特谱分析(Hilbert Spectral Analysis, HSA).本征模式分量必须满足下面两个条件:

(1)在整个数据序列中,极值点的数目与过零点的数目必须相等或至多相差一个;

(2)数据序列极大值点确定的上包络和极小值点所确定的下包络关于时间轴对称.

经验模式分解的主要过程如下:

(1)找到待分析信号的全部极大值和极小值点,利用三次样条函数分别把他们拟合为该信号的上下包络线,计算出两包络线的均值,进而求出待分析信号和均值的差值 h ;

(2)若 h 不满足IMF的要求,则重复上述过程若干次,使得新的 h 满足IMF的条件;若 h 满足IMF的要求,则令 h 为原信号的第1个IMF,并求出原信号与该IMF的差值 r ;

(3)将 r 作为待分解信号,重复以上过程,直到所剩余信号为单调信号为止.

从上面叙述可以知道EMD分解过程的时间复杂度为 $O(n)$,其中 n 为时间序列长度.

3 增量核空间独立向量组合预测算法

假设有时间序列 $\{x_t, t=1, 2, \dots\}$,通过时间窗长度为 l 的时间窗将时间序列分割成为如下形式的数据 $\{(\mathbf{x}_t, \mathbf{y}_t) | \mathbf{x}_t = (x_{t-l}, x_{t-l+1}, \dots, x_{t-1})^T, \mathbf{y}_t = x_t, t=l+1, \dots\}$,其中 \mathbf{x}_t 相当于输入变量, \mathbf{y}_t 相当于输出变量.

预测问题的标准方法可以使用形如 $f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle$ 简单参数化形式表示,其中 \mathbf{w} 为参数向量, ϕ 为映射函数, $\langle \cdot, \cdot \rangle$ 表示内积.函数 ϕ 可以为线性函数,也可以为非线性函数,函数 ϕ 将原始数据空间的数据映射到一个Hilbert空间或特征空间.当函数 ϕ 为线性函数时,预测问题等效于线性滤波器构造问题,如FIR、RLS算法等,但这些模型不适合于对非线性时间序列预测,所以为了对非线性时间序列进行有效预测应该取函数 ϕ 为非线性映射函数^[8,9].这样预测问题可以描述为最小化目标泛函 $\sum_{i=1}^l (y_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle)^2$ 求 \mathbf{w} .

求解如下优化问题:

$$\min L(\mathbf{w}) = \sum_{i=1}^l (y_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle)^2 = \|\Phi_l^T \mathbf{w} - \mathbf{y}_l\|^2 \quad (1)$$

可得最优参数 \mathbf{w}_l ,其中 $\mathbf{y}_l = (y_1, \dots, y_l)^T$, $\Phi_l = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_l)]$.由

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0 \quad (2)$$

$$\text{可得} \quad \mathbf{w}_l = (\Phi_l \Phi_l^T)^{-1} \Phi_l^T \mathbf{y}_l = \mathbf{K}_l^{-1} \Phi_l^T \mathbf{y}_l \quad (3)$$

直接求解上述方程会带来几个问题:

(1)当数据集很大时保存 \mathbf{K}_l 等将占用很大的内存空间;

(2)模型的阶数目等于时刻 t 时数据集中数据的个数容易造成过适;

(3) \mathbf{K}_l 条件数过大,这使得矩阵求逆不稳定.

为了克服上述问题,通过使用近似线性依赖条件(Approximate Linear Dependence, ALD)概念^[9],本文提出一种增量核空间独立向量组合预测算法(Incremental Independent Vector Combination Predicting algorithm in Kernel Space, IIVCPKS).该算法首先使用一些在核空间独立的输入变量对后续输入变量进行拟合,然后根据得到的拟合权值对独立变量对应的输出变量进行组合,组合结果作为后续输入变量所对应的输出值,即预测值.此外,如果后续输入变量不能被已有的独立变量很好的拟合,则将该独立变量作为一个独立变量加入到独立变量集(设独立变量集表示为 \mathbf{D}_l).下面详细介绍IIVCPKS算法的原理.

假设在 t 时刻有时间序列数据库 $\mathbf{DB}_{t-1} = \{(\mathbf{x}_i, \mathbf{y}_i), i=1, 2, \dots, t-1\}$ 和核空间独立变量集合 $\mathbf{D}_{t-1} = \{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i), i=1, 2, \dots, m_{t-1}, (\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i) \in \mathbf{DB}_{t-1}\}$,则新到的一个输入向量 \mathbf{x}_t 可以使用 \mathbf{D}_{t-1} 中独立变量进行拟合,则可得相应的最优拟合系数 $\mathbf{a} = (a_1, a_2, \dots, a_{m_{t-1}})^T$:

$$\begin{aligned} \mathbf{a} &= \arg \min_{\mathbf{a}} \left\| \sum_{i=1}^{m_{t-1}} a_i \phi(\tilde{\mathbf{x}}_i) - \phi(\mathbf{x}_t) \right\| \\ &= \arg \min_{\mathbf{a}} \left\{ \sum_{i,j=1}^{m_{t-1}} a_i a_j k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j) - 2 \sum_{j=1}^{m_{t-1}} a_j k(\mathbf{x}_j, \mathbf{x}_t) + k(\mathbf{x}_t, \mathbf{x}_t) \right\} \\ &= \arg \min_{\mathbf{a}} \{ \mathbf{a}^T \mathbf{K}_{t-1} \mathbf{a} - 2 \mathbf{a}^T \tilde{\mathbf{K}}_{t-1}(\mathbf{x}_t) + k_{tt} \} \end{aligned} \quad (4)$$

由式(4)可得 $\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{K}}_{t-1}(\mathbf{x}_t)$.为了使 $\tilde{\mathbf{K}}_{t-1}$ 不具有 \mathbf{K}_l 条件数太大、占用内存过多和易造成过拟合等缺点,应该限制 \mathbf{D}_t 的规模,在此引入近似线性依赖(Approximate Linear Dependence, ALD)条件:

$$\delta_t \stackrel{\text{def}}{=} \min_{\mathbf{a}} \left\| \sum_{j=1}^{m_{t-1}} a_j \phi(\tilde{\mathbf{x}}_j) - \phi(\mathbf{x}_t) \right\|^2 \leq v \quad (5)$$

由式(5)可得 $\delta_t = k_{tt} - \tilde{\mathbf{K}}_{t-1}^{-1}(\mathbf{x}_t)^T \mathbf{a}_t \leq v$.

如果 $\delta_t > v$,说明目前 \mathbf{D}_{t-1} 数据不能很好的表示 \mathbf{x}_t ,就是说 $\phi(\mathbf{x}_t)$ 与 $\phi(\tilde{\mathbf{x}}_i)$ 在高维Hilbert空间中线性独立,应该将 \mathbf{D}_{t-1} 更新为 $\mathbf{D}_t = \mathbf{D}_{t-1} \cup \{\mathbf{x}_t\}$,将 m_t 更新为 $m_t = m_{t-1} + 1$,从而为更好地表示下一个数据做准备.如果 $\delta_t > v$,说明目前 \mathbf{D}_{t-1} 数据能很好的表示 \mathbf{x}_t ,就是说 $\phi(\mathbf{x}_t)$ 与 $\phi(\tilde{\mathbf{x}}_i)$ 线性依赖,应该将 \mathbf{D}_{t-1} 更新为 $\mathbf{D}_t = \mathbf{D}_{t-1}$ 以为表示下一个数据做准备.对应于 \mathbf{D}_t 两种更新情况, $\tilde{\mathbf{K}}_t$ 的更新也可以分为两种情况:

(1)情况一:当 $\delta_t \leq v$ 时, $\mathbf{D}_t = \mathbf{D}_{t-1}$, $m_t = m_{t-1}$, $\tilde{\mathbf{K}}_t = \tilde{\mathbf{K}}_{t-1}$, $\tilde{\mathbf{K}}_t^{-1} = \tilde{\mathbf{K}}_{t-1}^{-1}$.

(2)情况二:当 $\delta_t > v$ 时, $\mathbf{D}_t = \mathbf{D}_{t-1} \cup \{\mathbf{x}_t\}$, $m_t = m_{t-1} + 1$,此时 $\tilde{\mathbf{K}}_t \neq \tilde{\mathbf{K}}_{t-1}$,则

$$\tilde{\mathbf{K}}_t = \begin{bmatrix} \tilde{\mathbf{K}}_{t-1} & \tilde{\mathbf{K}}_{t-1}(\mathbf{x}_t) \\ \tilde{\mathbf{K}}_{t-1}(\mathbf{x}_t)^T & k_{tt} \end{bmatrix}$$

$$\Rightarrow \tilde{\mathbf{K}}_t^{-1} = \frac{1}{\delta_t} \begin{bmatrix} \delta_t \tilde{\mathbf{K}}_{t-1}^{-1} + \mathbf{a}_t \mathbf{a}_t^T & -\mathbf{a}_t \\ -\mathbf{a}_t^T & 1 \end{bmatrix} \quad (6)$$

综上所述,增量核空间独立向量组合预测算法(Incremental Independent Vector Combination Predicting algorithm in Kernel Space, IIVCPKS)可描述如下:

算法: IIVCPKS 算法

输入: $v, \mathbf{DB} = \mathbf{D}_{t-1} \cup \{x_t\}, m_{t-1}, \mathbf{D}_{t-1}, \tilde{\mathbf{y}}_{t-1} = [\tilde{y}_1, \dots, \tilde{y}_{t-1}]^T, \tilde{\mathbf{K}}_{t-1}, \tilde{\mathbf{K}}_{t-1}^{-1}$

输出: $\hat{y}_t, m_t, \mathbf{D}_t, \tilde{\mathbf{y}}_t, \tilde{\mathbf{K}}_t, \tilde{\mathbf{K}}_t^{-1}$

步骤:

(1) 计算 $\tilde{\mathbf{K}}_{t-1}(x_t)$;

(2) 计算 $\mathbf{a}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{K}}_{t-1}(x_t)$ 和 $\delta_t = k_u - \tilde{\mathbf{K}}_{t-1}(x_t)^T \mathbf{a}_t$;

(3) 使用 $\hat{y}_t = \mathbf{a}_t^T \tilde{\mathbf{y}}_{t-1}$ 计算预测值;

(4) if $\delta_t \leq v$

$\mathbf{D}_t = \mathbf{D}_{t-1}, m_t = m_{t-1}, \tilde{\mathbf{K}}_t = \tilde{\mathbf{K}}_{t-1}, \tilde{\mathbf{K}}_t^{-1} = \tilde{\mathbf{K}}_{t-1}^{-1}$;

else

$\mathbf{D}_t = \mathbf{D}_{t-1} \cup \{x_t\}, m_t = m_{t-1} + 1$;

使用式(6)计算 $\tilde{\mathbf{K}}_t$ 和 $\tilde{\mathbf{K}}_t^{-1}$;

假设下一步观测或测量可得 x_t 的真实输出为 y_t , 则将 $\tilde{\mathbf{y}}_t = [\tilde{\mathbf{y}}_{t-1}^T, y_t]^T$;

end if

图 1 增量核空间独立向量组合预测算法描述

从图 1 可以知道 IIVCPKS 的时间复杂度为 $O(m^2)$, 其中 m 为 t 时刻各个数据在高维 Hilbert 空间中独立变量的个数。

4 复杂时间序列的新型预测模型

本文提出复杂时间序列的新型预测模型的工作流程如下。

(1) 使用经验模式分解对非线性非平稳时间序列进行 n 阶分解, 即分解出从高频到较低频的 n 个分量;

(2) 采用高斯函数作为核函数, 使用 IIVCPKS 算法对各个分量进行预测建模。由于经验模式分解为近似分频分解, 这样各个分量所需要拟合函数集的容量(即 VC 维)近似 $1/2$ 倍递减, 而核函数的参数决定了该容量, 所以各个分量所采用的最优核函数参数的对数取值应该呈现近似线性关系, 这样只需要进行两次交叉验证就可以求取各个分量所需的近似最优参数, 这样可以快速对不同的分量采用相应不同参数值的核函数进行 IIVCPKS 预测建模;

(3) 对各个分量的预测结果进行等权求和集成为综合预测结果。

因为 EMD 分解和 IIVCPKS 算法都具有很小的时间复杂度, 所以新预测模型非常适合于实时性要求比较高的应用场合。

5 实验仿真

为了验证新型预测模型的有效性, 采用了 4 个实际的非线性、非平稳随机时间序列进行实验, 并与单一分辨率下的 I

IVCPKS 预测模型及目前普遍使用且预测效果良好的 RBF 神经网络做对比实验。四个实际复杂时间序列为: Sunspot database(SD)^[10], Darwin Sea Level Pressure(DSLP)^[10], Poland Electric Demand time series(PED)^[10], Ehemet Packet 时间序列数据(EP)^[11]。

由于以太网数据包时间序列在这些非线性时间序列中最具典型的非线性性和非对称性等特点, 所以在此仅给出针对以太网数据包时间序列使用不同预测模型得到的预测结果图形显示, 如图 2 所示。图 2 中横坐标为时间轴, 纵轴为归一化的真实时间序列及各预测模型预测结果的幅值。图 2 的第一个子图给出了 RBF 神经网络和新模型的预测结果显示, 与 RBF 神经网络预测模型相比, 新模型的预测序列更逼近于真实序列, 可见新模型预测性能优于 RBF 神经网络预测模型。图 2 的第二个子图给出了新模型和不分解而直接采用 I-IVCPKS 模型的预测结果, 与单一 IIVCPKS 模型相比, 新模型的预测序列更逼近于真实序列, 可见新模型的预测性能优于单一 IIVCPKS 模型。

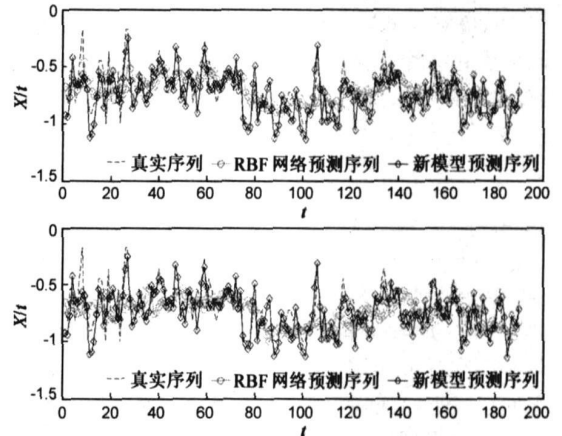


图 2 以太网数据包时间序列数据两种预测算法实验结果显示

为了全面考察某一预测方法的预测性能, 本文采用两种被普遍使用的性能评价标准: 平均绝对误差(Mean absolute error, MAE)和规范化根均方误差(Normalized root mean square error, NRMSE)。MAE 的定义式为:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x(i) - \hat{x}(i)| \quad (7)$$

其中 n 为预测集数据个数, $x(i)$ 为真实值, $\hat{x}(i)$ 为预测值。NRMSE 的定义式为:

$$\text{NRMSE} = \frac{\sqrt{\frac{1}{n} \sum_{i=1}^n [x(i) - \hat{x}(i)]^2}}{\sqrt{\frac{1}{n} \sum_{i=1}^n [x(i) - \bar{x}]^2}} \quad (8)$$

其中 n 为预测集数据个数, $x(i)$ 为真实值, $\hat{x}(i)$ 为预测值, \bar{x} 为序列均值。

四个混沌时间序列的预测误差分析结果如表 1 所示。由表 1 可知, 使用新模型进行预测的 MAE 或 NRMSE 明显小于直接采用 IIVCPKS 或 RBF 神经网络进行预测的 MAE 或 NRMSE。可见新预测模型的预测性能优于单一 IIVCPKS 或

RBF 神经网络预测模型的预测性能, 其原因在于在各个本征模式下的频率成分或波形变化与原始信号相比要简单, 从而更容易预测并且得到更好的预测结果.

表 1 不同非线性时间序列预测的 MAE 和 NRMSE 误差

数据集	MAE			NRMSE		
	RBF	IIVCPKS	EMD + IIVCPKS	RBF	IIVCPKS	EMD + IIVCPKS
EPS	0.1277	0.1432	0.0856	0.3702	0.3942	0.2366
PED	0.0376	0.0715	0.0253	0.0056	0.0090	0.0037
DSLP	0.1251	0.1763	0.1274	0.0445	0.0634	0.0456
SD	0.1350	0.1755	0.1017	0.0736	0.0901	0.0553

6 结论

本文提出的复杂时间序列的新型多分辨率增量预测模型, 实现了对非线性非平稳随机信号的分解, 并对各分量分别进行 IIVCPKS 预测建模. 使用 4 个实际的混沌时间序列进行验证其有效性的结果表明: 新型预测模型是一种适合于非线性、非平稳时间序列预测行之有效的预测方法, 预测性能要好于单一 IIVCPKS 及单一 RBF 神经网络, 因此具有广阔的应用前景. 但是 EMD 算法的端点效应在一定程度上会影响新型预测模型的预测性能, 所以对 EMD 算法端点效应对预测性能的影响及新的端点效应处理方法等的研究将是未来研究的方向.

参考文献:

- [1] Lendasse A, et al. Vector quantization: a weighted version for time series forecasting[J]. Future Generation Computer Systems, 2005, 21(7): 1056–1067.
- [2] Ljung L. System Identification Theory for User [M]. Prentice Hall, 1987. 1–115.
- [3] Versace M, et al. Predicting the exchange traded fund DIA with a combination of genetic algorithms and neural networks[J]. Expert Systems with Applications, 2004, 27(3): 417–425.
- [4] Mukherjee S, et al. Nonlinear prediction of chaotic time series using support vector machines[A]. Proceeding of the IEEE Workshop on Neural Networks for Signal Processing [C]. Amelz Island, 1997. 511–520.
- [5] Mitani Y, et al. Time series prediction of acoustic signals using neural network model and wavelet shrinkage[A]. Proceedings of the Tenth International Congress on Sound and Vibration [C]. Stockholm, Sweden: IIAV, 2003. 4189–4196.
- [6] Ince H, et al. Kernel principal component analysis and support vector machines for stock price prediction[A]. Proceeding of the IEEE International Joint Conference on Neural Networks [C]. Budapest, Hungary, 2004: 2053–2058.

- [7] Huang N E, et al. Applications of Hilbert Huang transform to nonstationary financial time series analysis [J]. Applied Stochastic Models in Business and Industry, 2003, 19(3): 245–268.
- [8] 解应春, 王海清, 李平. RKRLS 及在混炼胶质量建模与预测中的应用研究[J]. 浙江大学学报, 2004, 38(8): 941–945.
Xie Ying chun, Wang Hai qing, et al. RKRLS and its application to modeling and prediction of rubber compound quality[J]. Journal of Zhejiang University (Engineering Science), 2004, 38(8): 941–945. (in Chinese)
- [9] Engel Y, et al. The kernel recursive least squares algorithm[J]. IEEE Transaction on Signal Processing, 2004, 52(8): 2275–2285.
- [10] Time Series Prediction group[EB/OL]. <http://www.cis.hut.fi/projects/tsp/?page=Timeseries>, 2006 05 08.
- [11] Ethernet Packet[EB/OL]. http://math.bu.edu/people/murad/methods/time_series/index.html#Ethernet, 2006 05 08.

作者简介:



王 军 男, 1976 年生于浙江江山, 哈尔滨工业大学自动化测试与控制系博士生, 主要研究方向数据挖掘、信号与信息处理和智能故障诊断理论等.
E-mail: wangjunhit@chinaacc.com



彭喜元 男, 1961 年生于内蒙古四子王旗, 哈尔滨工业大学自动化测试与控制系教授、博士生导师, 主要研究方向为自动化测试技术和智能故障诊断理论等.



彭 宇 男, 1973 年生于陕西西安, 哈尔滨工业大学自动测试与控制系副教授, 主要研究领域为信号处理、计算智能和智能故障诊断理论.