

文章编号: 1001-0920(2013)04-0501-05

基于独立成分分析的含噪声时间序列预测

杨臻明¹, 岳继光¹, 王晓保², 萧蕴诗¹

(1. 同济大学 电子与信息工程学院, 上海 201804; 2. 上海申通轨道交通研究咨询有限公司, 上海 201103)

摘要: 提出一种基于独立成分分析(ICA)的最小二乘支持向量机(LS-SVM), 用于时间序列的多步超前独立预测. 用ICA估计预测变量中的独立成分(IC), 用不含噪声的IC重新构建时间序列. 利用 k -最近邻法(k -NN)减小训练集的规模, 提出一种新的距离函数以降低LS-SVM训练过程的计算复杂度, 并用约束条件对预测值进行后处理. 使用基于ICA的LS-SVM、普通LS-SVM与反向传播神经网络(BP-ANN), 对多个时间序列进行对比预测实验. 实验结果表明, 基于ICA的LS-SVM的预测性能优于普通LS-SVM和BP-ANN.

关键词: 独立成分分析; 时间序列预测; k -最近邻法; 最小二乘支持向量机

中图分类号: TP273

文献标志码: A

Noisy time series prediction using independent component analysis

YANG Zhen-ming¹, YUE Ji-guang¹, WANG Xiao-bao², XIAO Yun-shi¹

(1. College of Electronic and Information Engineering, Tongji University, Shanghai 201804, China; 2. Shanghai Shentong Rail Transit Research and Consultancy Co Ltd, Shanghai 201103, China. Correspondent: YANG Zhen-ming, E-mail: yangzhenmingtk@yahoo.com.cn)

Abstract: A least square support vector machine (LS-SVM) based on the independent component analysis(ICA) is proposed to predict noisy non-stationary time series. ICA is used to estimate the independent components(IC) in the forecasting variables. After identifying and removing the ICs containing the noise, the rest of the ICs are then used to reconstruct the forecasting variables which contain less noise. A k -nearest neighbors(k -NN) approach is used to reduce the size of training dataset and a new distance function is defined. By selecting similar instances in the training dataset, the complexity of training a LS-SVM is reduced significantly. A boundary constraint component is developed to limit the predicted values to a reasonable range. The experimental results show that the proposed approach outperforms both traditional LS-SVM and BP-artificial neural network(BP-ANN) in the prediction performance of several time series.

Key words: independent component analysis; time series prediction; k -nearest neighbors; least square support vector machine

0 引言

时间序列预测在过程控制、金融证券、交通流量、电力负荷^[1-2]等许多领域得到了广泛应用. 实际的时间序列大都表现出非线性、非平稳特征, 并且包含了大量噪声^[3]. 由支持向量机(SVM)发展而来的最小二乘支持向量机(LS-SVM)^[4]因其优越的性能已成为解决时间序列预测问题的重要手段^[5]. LS-SVM结合人工神经网络(ANN)提高了时间序列的预测性能^[6], 降低了预测器的计算复杂度, 同时也增加了选择输入变量的难度. 噪声是影响时间序列LS-SVM建模的关键问题之一, 忽略噪声影响将导致模型出现过拟合或

欠拟合现象, 降低模型的泛化能力^[7]. 有关筛选训练集与降低序列噪声等方面的研究工作目前还处于起步阶段.

本文提出一种基于LS-SVM的时间序列直接预测方法. 运用独立成分分析(ICA)检测并降低时间序列的噪声, 提高LS-SVM的预测性能. 该方法运用ICA估计预测变量中的独立成分(IC), 在识别与剔除含噪声的IC后, 用不含噪声的IC重新构建时间序列. 然后运用 k -最近邻法(k -NN)选择最近邻样本来构成LS-SVM的简化训练集, 其距离函数结合了欧氏距离和时间序列的一阶差分. 运用基于 k -最近邻筛选

收稿日期: 2011-12-11; 修回日期: 2012-02-21.

基金项目: 国家自然科学基金项目(40872090).

作者简介: 杨臻明(1982-), 男, 博士生, 从事非线性系统、神经网络的研究; 岳继光(1961-), 男, 教授, 博士生导师, 从事过程控制、计算机控制等研究.

和独立成分分析的最小二乘支持向量机 (k -NN-ICA-LS-SVM), 对多个时间序列进行了预测实验, 并比较了该方法与其他预测方法之间的差异.

1 滑窗与训练矩阵

时间序列是由特定时间点 t 上的观测值 x_t 组成的有序队列. 一个长度为 t 的时间序列可以表示为 $X = [x_1, x_2, \dots, x_t]$, 并将其中的一段 $[x_{t-p}, x_{t-p+1}, \dots, x_t]$ 记为 X_{t-p}^t . h 步超前预测即是在已知过去 p 个观测值 X_{t-p+1}^t 的情况下, 对序列后续 h 个值 X_{t+1}^{t+h} 进行预测. 区别于递归方法, 独立预测方法的每一预测步长都对应于独立预测模型^[8], 第 i 步的预测值可由预测模型 f_i 表示为

$$\hat{x}_{t+i} = f_i(x_{t-p+1}, \dots, x_{t-1}, x_t), 1 \leq i \leq h. \quad (1)$$

对于时间序列 X_t , 可以利用如图1所示的滑窗来建立训练集 D , 每个滑窗对应于训练集中的一个样本, 其结果相当于将 1 维时间序列展开成 2 维矩阵. 假设滑窗的长度为 $p+h$, 时间序列长度为 T , 则通过滑窗法获得的训练集 D 是一个 $(T-p-h+1) \times (p+h)$ 矩阵. 矩阵 D 的前 p 列代表所有预测模型的训练输入值, 第 $p+i$ 列代表第 i 步预测模型 f_i 的训练输出值. 矩阵 D 的每一行代表一个训练样本, 其前 p 个值代表输入, 后 h 个值分别是 h 个预测模型的输出.

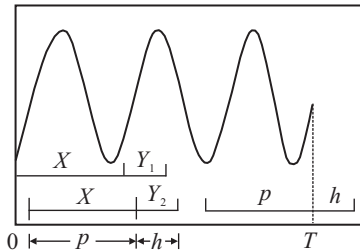


图 1 滑窗与训练集

2 独立成分分析

ICA 是一种统计信号处理方法, 可以在缺乏先验知识的条件下, 仅通过包含未知信号的观测数据来恢复源信号^[9]. 本文利用 ICA 消除时间序列包含的噪声. 设 X 是通过滑窗法 (见第 1 节) 获得的一个 $(T-p-h+1) \times (p+h)$ 训练集矩阵, 令 $m = T-p-h+1$, $n = p+h$, 则 X 可以表示为 $m \times n$ 矩阵 $[x_1, x_2, \dots, x_m]^T$, 其中每个行向量 x_i 代表一个长度为 n 的训练样本. 在 ICA 中, 矩阵 X 可以表示为^[9]

$$X = AS = \sum_{i=1}^m a_i s_i, \quad (2)$$

其中: A 是 $m \times m$ 混合矩阵, S 是 $m \times n$ 源矩阵. 行向量 s_i 表示无法从混合信号 x_i 中直接观测到的隐含源信号. ICA 的目标是找到 $m \times m$ 逆混合矩阵 W 以恢复源矩阵 S , 即

$$Y = [y_i] = WX, \quad (3)$$

其中 y_i 是矩阵 Y 的行向量. 为了用 y_i 估计独立的隐含源信号 x_i , 在统计意义下 y_i 之间必须相互独立, 因此将其称为独立成分 (IC). 当逆混合矩阵 W 是混合矩阵 A 的逆矩阵时, 即当 $W = A^{-1}$ 时, 可以将 IC (y_i) 作为隐含源信号 s_i 的估计值.

为识别包含噪声的 IC, 本文采用相对 Hamming 距离 (RHD)^[10] 作为评价 IC 的指标. RHD 的定义为

$$\text{RHD} = \frac{1}{n-1} \sum_{i=1}^{n-1} [R_i(t) - \hat{R}_i(t)]^2. \quad (4)$$

其中: $R_i = \text{sgn}[T_i(t+1) - T_i(t)]$, $\hat{R}_i = \text{sgn}[A_i(t+1) - A_i(t)]$, $\text{sgn}(r)$ 是符号函数, T_i 是时间序列的实际值, A_i 是预测值, n 是时间序列长度. RHD 可以用来评价时间序列之间的相似度.

假设已求得 $m \times m$ 逆混合矩阵 W , 则根据式 (3) 可得 m 组 IC. 用除 y_k 以外的 $m-1$ 组 IC 重构信号矩阵, 并对所有 m 组 IC 都进行相同计算, 有

$$\hat{X} = \sum_{i=1, i \neq k}^m a_i y_i, 1 \leq k \leq m. \quad (5)$$

其中: $\hat{X} = [\hat{x}_1, \dots, \hat{x}_{k-1}, \hat{x}_{k+1}, \dots, \hat{x}_m]^T$ 是重构的 $(m-1) \times n$ 信号矩阵, \hat{x}_i 是重构行向量, a_i 是混合矩阵 A 的第 i 列, $A = W^{-1}$. 用 RHD 值评价原始行向量 x_i 与对应的重构行向量 \hat{x}_i 之间的相似度, RHD 值越小, 说明 x_i 与 \hat{x}_i 的主要特征越相似, 同时也说明用于重构时间序列的 IC 包含原始序列更多的特征. 换言之, 排除在重构 IC 之外的 IC (即 y_k 对应的 IC) 包含更多的噪声. 因此剔除该 IC 后, 用最小 RHD 值对应的 IC 重构时间序列, 即可降低时间序列所包含的噪声.

如果时间序列的信噪比较低, 则剔除 1 组 IC 后重构时间序列的 RHD 值将比较平均且相对较大, 不易从中识别包含噪声的 IC. 这时可将重构时间序列的 IC 组数减少至 $m-2$, 通过遍历所有 C_m^{m-2} 个 IC 组合并重复上述步骤计算 RHD, 直至 RHD 小于预设值或触发停止条件^[10].

为了求解逆混合矩阵 W , 需要运用一些优化方法. 若将 IC 的独立性度量作为目标函数, 则 ICA 建模可以归结为求解一个最优化问题, 目前已有一些这方面的算法^[11-12]. IC 的统计独立性可由非正态分布描述^[9], 通过负熵度量, 即

$$J(y) = H(y_g) - H(y), \quad (6)$$

其中 y_g 是一正态随机向量, 其协方差矩阵与 y 相同. 随机向量的概率密度为 $p(y)$, H 是 y 的熵, 其定义为

$$H(y) = - \int p(y) \log p(y) dy. \quad (7)$$

该熵是非负的, 当且仅当 y 服从正态分布时等于 0. 由于负熵的计算非常困难, 通常采用如下的近似值^[9]:

$$J(y) \approx [E\{G(y)\} - E\{G(v)\}]^2. \quad (8)$$

其中: $G(y) = \exp(-y^2/2)$, v 是一个平均值为 0、方差为 1 的标准正态变量, y 是一个平均值为 0、方差为 1 的任意随机变量. 本文采用一种快速 ICA 算法^[11] 计算逆混合矩阵 W .

在 ICA 建模之前有 2 个预处理步骤^[9]. 首先, 将输入矩阵 X 的每一行 x_i 减去该行的平均值, 即 $x_i \leftarrow (x_i - E(x_i))$. 其次, 消除输入矩阵 X 的二阶统计量, 即

$$Z = 2C_x^{-1}X. \quad (9)$$

其中: $C_x = \sqrt{E(xx^T)}$ 是输入向量协方差矩阵的平方根, C_x^{-1} 是 C_x 的逆矩阵. 经过预处理得到的输入矩阵 Z , 其行向量 z 之间不相关, 且具有单位方差, 即

$$E(zz^T) = I. \quad (10)$$

本文研究的时间序列都经过上述预处理.

3 k -NN 与混合距离测度

通常 k -NN 应用于时间序列预测时, 先对输入 x 确定一个邻域范围, x 对应的输出值 y 取该邻域中所包含的 k 个样本的输出平均值. 假定相似输入输出之间的映射关系也是相似的, 则 k -NN 选取最接近测试样本的 k 个训练样本构成简化训练集, 该训练集用于训练时间序列预测模型.

通常使用欧氏距离度量样本之间的相似性. 然而许多实际的时间序列具有非平稳特性^[3], 在其演化过程中不具有固定不变的均值, 同时, 序列的某一部分与其他部分又非常相似, 这时欧氏距离便无法准确刻画样本之间的相似性. 通过对时间序列作适当的差分可使之平稳化, 然后再度量样本之间的相似度, 从而消除了时间序列变化趋势的影响^[13]. 考虑到计算复杂度, 本文采用一阶差分方法. 假设有从 t 到 $t+T$ 的训练矩阵 D , 对应矩阵的第 j 行为 $(x_{t+j}, x_{t+j+1}, \dots, x_{t+j+p-1}, x_{t+j+p}, \dots, x_{t+j+p+h-1})$, 其中 $1 \leq j \leq (T-p-h+1)$. 与其对应的输入向量为 $(x_{t+j}, x_{t+j+1}, \dots, x_{t+j+p-1})$, 输入向量的一阶差分为

$$(d_{t+j}, \dots, d_{t+j+p-2}) = (x_{t+j+1} - x_{t+j}, \dots, x_{t+j+p-1} - x_{t+j+p-2}),$$

该向量是一个 $p-1$ 维行向量.

对于 $T+1$ 处的 p 维输入向量 $(x_{T+1}, \dots, x_{T+p-1}, x_{T+p})$, 其与训练集中每个样本的欧氏距离可以表示为

$$E(j) = \sqrt{\sum_{i=1}^p (x_{T+i} - x_{t+j+i-1})^2}. \quad (11)$$

测试向量的一阶差分为 $(d_{T+1}, \dots, d_{T+p-1}) = (x_{T+2} - x_{T+1}, \dots, x_{T+p} - x_{T+p-1})$, 其同样也是一个 $p-1$ 维

行向量. 其与训练集中每个样本一阶差分的欧氏距离可以表示为

$$D(j) = \sqrt{\sum_{i=1}^{p-1} (d_{T+i} - d_{t+j+i-1})^2}. \quad (12)$$

E 和 D 中的向量个数都为 $(T-p-h+1)$, 归一化 E 和 D 的组合 D^* 构成了 k -NN 的混合距离测度, 即

$$D^*(j) = \frac{E(j) - E_{\min}}{E_{\max} - E_{\min}} + \frac{D(j) - D_{\min}}{D_{\max} - D_{\min}}, \quad (13)$$

其中 E_{\min} 、 E_{\max} 、 D_{\min} 、 D_{\max} 分别为 E 和 D 的最小值和最大值. 混合距离测度考虑了序列非平稳特征的影响, 距离测度最小的 k 个样本构成了 LS-SVM 的简化训练集.

4 基于 k -NN 的 LS-SVM

本文利用 LS-SVM 建立时间序列的非线性回归模型^[5]. 考虑初始权值空间中有如下形式的模型:

$$y(x) = w^T \varphi(x) + b. \quad (14)$$

其中: $x \in \mathbf{R}^p$, $y \in \mathbf{R}$, $\varphi(x)$ 是将输入 x 变换到高维特征空间的非线性映射函数. 给定训练集 $\{x_j, y_j\}_{j=1}^N$, 最优化问题可以表示为

$$\min_{w, b, e} J_P(w, e) = \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{j=1}^N e_j^2; \quad (15)$$

$$\text{s.t. } y_j = w^T \varphi(x_j) + b + e_j, \quad j = 1, 2, \dots, N. \quad (16)$$

其中 e_j 是误差变量. 当 w 的维数趋于无穷时, 需要建立并解决上述问题的 Lagrange 对偶问题. 该最优化问题的解是

$$y(x) = \sum_{j=1}^N \alpha_j K(x, x_j) + b, \quad (17)$$

其中 $K(x, x_j)$ 是核函数, 其定义为

$$K(x_j, x_l) = \varphi(x_j)^T \varphi(x_l). \quad (18)$$

本文采用高斯核函数, 即

$$K(x, x_j) = \exp\left(-\frac{\|x - x_j\|^2}{2\sigma^2}\right). \quad (19)$$

高斯核函数有两个可调参数 γ 和 σ , 其数值通常在实验中确定.

本文采用独立预测方法, 对每一个测试样本建立独立的预测模型. 同时, 选择最接近测试样本的 k 个训练样本, 简化 LS-SVM 的训练集, 有效降低了建立 LS-SVM 回归模型的复杂度.

建立 k -NN-ICA-LS-SVM 共有 5 个步骤: 1) 用 ICA 重构时间序列, 降低数据所含噪声; 2) 对于测试集 T 中的每个样本, 以式 (13) 定义的混合距离测度在训练集 D 中选择最接近测试样本的 k 个样本, 其中训练集 D 由该测试样本的前 T 个数据点组成; 3) 用该测试样本的简化训练集 D' 建立每个预测步长的 LS-SVM 回归模型, h 步超前预测对应于建立 h 个独

立回归模型; 4) 将测试样本输入 h 个回归模型, 得到 h 步预测值; 5) 应用约束条件验证预测值.

当预测步长 h 较大时, 预测值可能会超出合理范围, 因此需要用约束条件来处理 LS-SVM 模型产生的预测值. 将约束条件的上限 B_{up} 和下限 B_{low} 分别设置为

$$B_{up} = d_{max} + 0.02d_{std}, \quad (20)$$

$$B_{low} = d_{min} - 0.02d_{std}. \quad (21)$$

其中: d_{max} 和 d_{min} 分别代表训练样本的最大值和最小值, d_{std} 代表训练样本的标准差. 如果预测值大于约束条件的上限或小于下限, 则将该值设为新的上限或下限; 否则, 将其作为最终预测值.

5 时间序列预测实验

为了评价 k -NN-ICA-LS-SVM 方法的时间序列预测性能, 进行了多个预测实验, 并用该方法与普通 LS-SVM 和反向传播神经网络 (BP-ANN) 进行了对比. 实验采用 Matlab 作为计算平台, 版本为 R2010b.

5.1 实验数据

实验采用两类数据集, 包括 Mackey-Glass 时间序列以及 NNGC1 竞赛提供的 4 个时间序列. Mackey-Glass 时间序列^[14]由如下形式的延时微分方程产生:

$$\frac{dx(t)}{dt} = \frac{ax(t-\tau)}{1+x(t-\tau)^{10}} - bx(t), \quad (22)$$

该序列常被用来评价各种预测方法. 用 4 阶 Runge-Kutta 方法产生长度为 2 201 的人工时间序列, 初始值 $x(0) = 1.2$, $a = 0.2$, $b = 0.1$, $\tau = 17$, 其中后 2 000 个数据点用于本次实验. NNGC1 竞赛^[15]提供了许多不同类型的异方差非平稳时间序列, 本次实验采用了其中 4 个有关交通流量的数据集, 采样周期为 1 h, 每个时间序列的长度为 1 742.

5.2 误差检测

本文采用两个指标来评价时间序列的预测误差. 用 X 代表测试时间序列的 h 个真实值, \hat{X} 代表 h 步超前预测得到的 h 个预测值. 第 1 个误差指标是均方根误差 (RMSE), 即

$$e_{RMSE} = \sqrt{\sum_{t=1}^h (x_t - \hat{x}_t)^2 / h}. \quad (23)$$

第 2 个误差指标是对称绝对值均差 (SMAPE), 其基于如下定义的相对误差:

$$e_{SMAPE} = \frac{1}{h} \sum_{t=1}^h \frac{|x_t - \hat{x}_t|}{(x_t + \hat{x}_t)/2}. \quad (24)$$

5.3 实验结果

如第 4 节所述, 在建立 k -NN-ICA-LS-SVM 的过程中, 训练序列长度 T , 输入向量长度 p , k -NN 参数 k ,

核函数参数 γ 与 σ , 预测步长 h 等一系列参数需要确定. 将比较实验的预测步长 h 设为 20, 各参数取值范围见表 1.

表 1 参数取值范围

参数	T	p	k	γ	σ
取值范围	[500, 1 000]	[3, 30]	[50, 150]	[0.02, 100]	[0.02, 50]

将 RMSE 最小时所对应的参数值作为预设值. 为了比较本文方法与普通 LS-SVM, 对同一时间序列两种方法采用相同的参数. 实验中各时间序列的参数选择情况见表 2.

表 2 参数预设值

时间序列	T	p	k	γ	σ
Mackey-Glass	700	25	80	30	50
NNGC1-1	600	20	70	10	10
NNGC1-2	600	20	60	5	10
NNGC1-3	600	25	110	10	10
NNGC1-4	600	20	70	5	10

采用 3 种方法对时间序列 NNGC1-4 进行 20 步超前预测, 预测结果见图 2. 图 2 显示了 1 123 个测试数据中前 100 个数据点的预测值与对应的真实值. 可以看出, k -NN-ICA-LS-SVM 方法的预测结果非常贴近真实数据, 优于普通 LS-SVM 方法与 BP-ANN 方法.

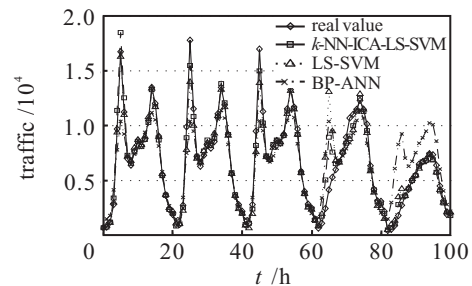


图 2 NNGC1-4 预测结果

从图 2 可以看出, BP-ANN 方法未能成功预测 NNGC1-4 序列随时间变化的趋势, 尤其是后 20 个数据点的预测误差较大. 普通 LS-SVM 方法虽然也表现出较好的预测性能, 但是由于训练集规模较大需要更多训练时间, 不适用于实时性要求较高的应用领域. 在本次实验中, 普通 LS-SVM 方法平均耗时 28.76 s 进行一次 20 步超前预测, 而 k -NN-ICA-LS-SVM 方法进行同样预测仅耗时 0.51 s.

对所有 5 个时间序列都进行了步长为 20 的超前预测实验, 预测结果的 RMSE 列于表 3, SMAPE 列于表 4. 实验结果表明, 对于所有时间序列, k -NN-ICA-LS-SVM 的预测误差都是最小的. SMAPE 指标由相对误差导出, 因此更具可比性. 人工时间序列 Mackey-Glass 不包含任何噪声, 其预测误差显著小于来自于实际问题的 NNGC1 序列.

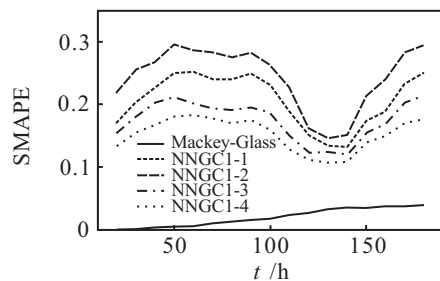
表3 RMSE 误差指标

时间序列	BP-ANN	普通 LS-SVM	k -NN-ICA-LS-SVM
Mackey-Glass	0.070 3	0.008 7	0.001 6
NNGC1-1	6594.2	4039.9	3608.4
NNGC1-2	155.84	113.1	104.52
NNGC1-3	7828.1	4771.5	4446.8
NNGC1-4	2396.6	1683	1491

表4 SMAPE 误差指标

时间序列	BP-ANN	普通 LS-SVM	k -NN-ICA-LS-SVM
Mackey-Glass	0.069 8	0.007 7	0.001 3
NNGC1-1	0.432 3	0.197 8	0.168 6
NNGC1-2	0.391 3	0.276 4	0.218 8
NNGC1-3	0.390 9	0.167 7	0.154 0
NNGC1-4	0.293 6	0.160 4	0.134 7

在进一步的实验中,将预测步长增加至180,以考察 k -NN-ICA-LS-SVM 的长期预测性能,预测结果的 SMAPE 见图3.可以看出,所有序列的预测误差都具有稳定性,不存在突变的情况.实验结果表明, k -NN-ICA-LS-SVM 在时间序列长期预测方面同样具有良好的性能.

图3 k -NN-ICA-LS-SVM 长期预测性能

6 结 论

针对时间序列预测在实际应用中遇到的问题,本文提出了一种 k -NN-ICA-LS-SVM 时间序列多步超前预测方法.采用 ICA 对训练集进行预处理,将数据集分解为 IC 的组合,用不含噪声的 IC 重构训练集,以降低数据集包含的噪声.在 k -NN 中使用一种混合距离测度对训练数据进行筛选,减小了训练集的规模.用约束条件对预测值进行后处理,以确保预测结果在具有实际物理意义的合理范围内.对多个时间序列进行了对比预测实验,实验结果表明,使用该方法可以获得更小的训练集,减少计算过程所需时间,并且预测误差显著低于普通 LS-SVM 方法和 BP-ANN 方法. k -NN-ICA-LS-SVM 在含噪声非平稳时间序列预测方面表现出良好的应用前景.

参考文献(References)

[1] Araujo R, Ferreira T A E. An intelligent hybrid morphological-rank-linear method for financial time series prediction[J]. Neurocomputing, 2009, 72(10/11/12): 2507-2524.

[2] 曹玉苹,田学民.基于SVM和Kalman预测的非线性系统故障预报[J].控制与决策,2009,24(3):477-480.
(Cao Y P, Tian X M. Nonlinear system fault prognosis based on SVM and Kalman predictor[J]. Control and Decision, 2009, 24(3): 477-480.)

[3] Yang H Q, Huang K Z, King I, et al. Localized support vector regression for time series prediction[J]. Neurocomputing, 2009, 72(10/11/12): 2659-2669.

[4] Adankon M M, Cheriet M, Biem A. Semisupervised least squares support vector machine[J]. IEEE Trans on Neural Networks, 2009, 20(12): 1858-1870.

[5] Sorjamaa A, Hao J, Reyhani N, et al. Methodology for long-term prediction of time series[J]. Neurocomputing, 2007, 70(16/17/18): 2861-2869.

[6] 李松,刘力军,解永乐.遗传算法优化BP神经网络的短时交通流混沌预测[J].控制与决策,2011,26(10):1581-1585.
(Li S, Liu L J, Xie Y L. Chaotic prediction for short-term traffic flow of optimized BP neural network based on genetic algorithm[J]. Control and Decision, 2011, 26(10): 1581-1585.)

[7] Cao L J. Support vector machines experts for time series forecasting[J]. Neurocomputing, 2003, 66(1/2/3): 321-339.

[8] Sapankevych N I, Sankar R. Time series prediction using support vector machines: a survey[J]. IEEE Computational Intelligence Magazine, 2009, 4(2): 24-38.

[9] Li X L, Adali T. Independent component analysis by entropy bound minimization[J]. IEEE Trans on Signal Processing, 2010, 58(10): 5151-5164.

[10] Cheung Y M, Xu L. Independent component ordering in ICA time series analysis[J]. Neurocomputing, 2001, 64(1/2/3): 145-152.

[11] Gao Q X, Zhang L, Zhang D, et al. Independent components extraction from image matrix[J]. Pattern Recognition Letters, 2010, 31(3): 171-178.

[12] Yu S N, Chou K T. Selection of significant independent components for ECG beat classification[J]. Expert Systems with Applications, 2009, 36(2): 2088-2096.

[13] Box G E P, Jenkins G M, Reinsel G C. Time series analysis: Forecasting and control[M]. 4th ed. New York: John Wiley & Sons, 2008: 109-116.

[14] Mackey M, Glass L. Oscillation and chaos in physiological control systems[J]. Science, 1977, 197(4300): 287-289.

[15] Crone S F. Artificial neural network & computational intelligence forecasting competition[EB/OL]. (2010-2-12) [2010-9-18]. <http://www.neural-forecasting-competition.com>.