

支持向量回归的颅内压时间系列无损估计方法

吴少智^{1,4}, 吴跃², 徐鹏³, 胡晓⁴

(1. 电子科技大学无锡研究院 成都 611731; 2. 电子科技大学计算机科学与工程学院 成都 611731;
3. 电子科技大学神经信息教育部重点实验室 成都 610054; 4. 加州大学洛杉矶分校 美国 CA 90025)

【摘要】在对时间序列数据挖掘框架进行研究时发现: 在利用线性映射函数刻画误差和特征间的关系时, 不能获得对颅内压力信号的精确估计。为了提高对颅内压估计的精确性, 本文采用支持向量回归构建存在于特征和误差间的非线性映射函数, 实验结果表明: 基于支持向量回归的非线性映射函数预测效果明显优于先前所采用的线性最小二乘法所构成的线性映射函数策略。

关键词 数据挖掘框架; 最小二乘法; 非线性映射函数; 支持向量回归

中图分类号 TP301.6

文献标识码 A

doi:10.3969/j.issn.1001-0548.2011.05.029

Support Vector Regression Based Time Series Mining Approach for Non-Invasive ICP Assessment

WU Shao-zhi^{1,4}, WU Yue², XU Peng³, and HU Xiao⁴

(1. Institute of Wuxi, University of Electronic Science and Technology of China Chengdu 611731;

2. School of Computer Science and Engineering, University of Electronic Science and Technology of China Chengdu 611731;

3. Key Laboratory for Neuro Information of Ministry of Education, University of Electronic Science and Technology of China Chengdu 610054;

4. Los Angeles, University of California USA CA 90025)

Abstract For the data mining based on time series estimation, the existed studies reveal that the Intra-Cranial Pressure (ICP) time series cannot be well estimated when the linear mapping function is used to delineate the relationship between error and feature. To improve the accuracy for ICP estimation, the non-linear support vector regression (SVR) is used to construct the nonlinear function between feature and error. The experiment results showed that the SVR based mapping function is superior to the linear least square based one.

Key words data mining framework; linear least squares; nonlinear mapping function; support vector regression

针对颅内压力信号的无损估计问题^[1-3], 曾经提出过一个数据挖掘框架^[4]用于对颅内压(Intra-Cranial Pressure, ICP)进行无损估计。首先从数据库中提取数据, 建立相应的训练数据集, 在构建相应的时间序列估计模型时, 采用系统辨识技术和特征抽取技术^[4], 从动脉血压(arterial blood pressure, ABP)和脑血流速度(cerebral blood flow velocity, CBFV)时间序列中提取必要的特征向量, 所提取的特征向量通过映射函数被映射到误差空间, 利用已有的各个估计模型得出对应ICP的估计误差, 然后根据所得出的估计误差, 结合从已有的模型数据库中选择相应的估计模型, 对ICP进行估计。由于误差估计的准确度直接影响着对颅内压估计的精度, 因此, 对特征和误

差间关系的描述, 即映射函数的构建, 是该框架中的一个关键技术。文献[5]首先研究了基于线性最小二乘法(linear least squares, LLS)构造的映射函数, 但该映射函数没有考虑真实信号可能受到噪声和伪信号的干扰。该文献还分别采用基于总体最小二乘法(total least squares, TLS)^[6]和奇异值截断分解(truncated singular value decomposition, TSVD)^[7]构造映射函数, 以消除可能的噪声影响, 从而提高对无损颅内压估计的精度。LLS、TLS、STR和TSVD这4种方法在本质上都属于线性范畴, 均采用简单的线性关系来刻画误差和特征间的关系, 线性函数的优点是方法简单、易于描述和表征。但是特征和误差间存在着复杂的关系, 超出了线性方法的学习能

收稿日期: 2010-03-03; 修回日期: 2010-07-19

基金项目: 国家863计划(2007AA01Z443); 四川省基础应用研究项目(2010JY0001)

作者简介: 吴少智(1972-), 男, 博士生, 主要从事数据挖掘、时间序列分析及生物医学信息分析、计算机网络等方面的研究。

力, 有必要采用更为复杂的非线性方法^[8]才能更好地刻画。

有多种方法可以建立非线性映射函数, 支持向量机方法(support vector machine, SVM)是其中常用的一种方法。SVM有别于传统方法, 其基于结构风险最小化原则, 最小化的是推广误差的上界, 通过在模型的复杂性和训练误差之间寻求平衡点, 得到一个最优的网络结构, 可较好地解决过学习和欠学习问题。另外, SVM的训练等价于一个线性约束二次规划问题, 从而可以得到唯一的全局最优解。正是由于上述优点, SVM具有良好的性能, 在解决各种非线性学习问题时取得了很大成功, 还被扩展应用到了回归、时间序列预测、故障诊断及生理信号处理等多个方面。

在本文中, 采用支持向量回归(support vector regression, SVR)分析的方法^[9-10]构建非线性映射函数, 在特征值和差异值之间建立更精确的关系, 并比较讨论线性和非线性映射函数对无损ICP估计的影响, 期望能够提高对颅内压的无损估计的精确度。

1 方法

1.1 数据挖掘框架简介

在该数据挖掘框架中建立相应的时间序列数据库, 假设有n个数据构成的时间序列对作为训练用的输入, 每个输入数据对由所期望的时间序列(desired time series, DTS)和相关时间序列(related time series, RTS)数据集构成。本文中, DTS为ICP时间序列, RTS为ABP和CBFV时间序列。数据挖掘框架从逻辑上被划分成在线仿真处理和离线学习两个部分。基于数据挖掘的时间序列预测框架的离线学习训练过程如

图1所示。该训练过程包括模拟仿真的2个过程: 1) 基于系统辨识建立时间序列预测模型, 选取合适的输入对; 2) 输入RTS, 用于模拟预测DTS。图1中, 在对第kth个DTS-RTS对的训练过程中, 将每一个输入对和相应的映射函数相关联。完成该训练过程需要重复n次处理, 以RTS分析中所获得的特征向量作为映射函数的输入, 把特征向量映射到误差空间, 训练结束后, 生成n个映射函数。在线性假设下, 对某一估计模型估计的误差向量和特征矩阵间的映射关系可以用下面的线性方程进行描述:

$$F \times b = e \tag{1}$$

式中, e 为误差向量; $F = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_M \end{bmatrix}$ e 为特征向量矩阵;

b 是与血液动力学相关的向量。在训练过程中, 通过已知的F和e对b进行估计。ICP估计值的精确与否, 与映射函数估计值的精度有直接关系。

在图1所示的训练过程中, 将同时测得的DTS和RTS组成对进入数据库条目, 训练过程显示所获得的映射函数针对的第K个数据条目。基于数据挖掘的时间序列预测框架的在线仿真处理流程如图2所示。图1中映射函数的训练结果作为图2中在线训练的映射函数。从RTS输入对(ABP和CBFV时间序列所构成的输入对)中抽取相应的特征向量, 然后通过已经训练好的映射函数, 将该特征向量映射到误差空间, 产生n个估计误差。基于估计误差, 采用最小差异准则策略选取, 选取训练数据库中最佳DTS-RTS对应的系统辨识模型, 对ICP进行估计。

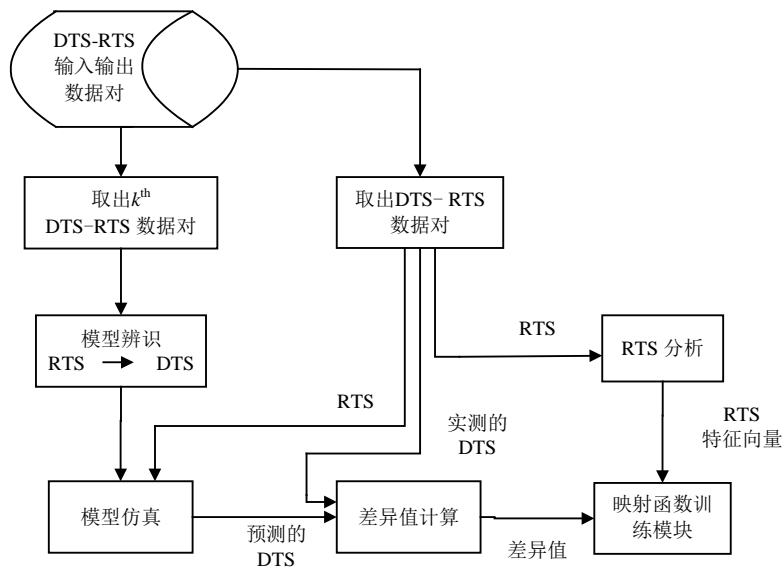


图1 基于数据挖掘的时间序列预测框架的训练过程

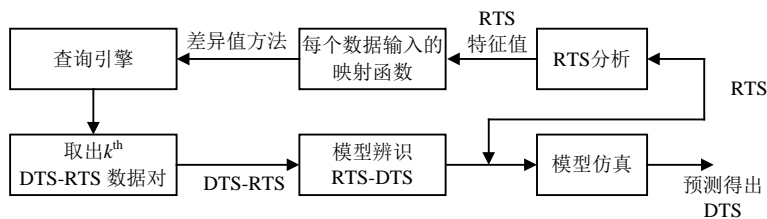


图2 基于数据挖掘的时间序列预测框架的在线仿真处理流程

1.2 误差函数

在时间序列预测框架中，根据预测的侧重点不同，可采用不同的误差度量方法。共有5种不同的差值计算方法，其中相似性误差测度方法是针对当前研究工作的效果较好的方法，其定义为：

$$e = 1 - \text{corr}(y^N, \hat{y}^N) \quad (2)$$

式中， $e \in [0, 2]$ ； y^N 代表规范化后的ICP序列； \hat{y}^N 则是 y^N 的相应的预测序列； $\text{corr}(x, y)$ 是时间序列 x 和 y 之间的相关系数。误差测度方法主要关注估计ICP和真实ICP间的匹配程度。

文献[5]采用如下公式计算输入信号 x 的值：

$$x^N = \frac{x - \bar{x}}{\bar{x}} \quad (3)$$

式中， \bar{x} 是输入信号 x 的平均值。式(3)中，由于某些输入时间序列的平均值可能存在为零的趋势，因而可能会造成式(3)的分母值非常小，可能趋近于零甚至等于零，造成 x^N 趋于无穷大，使输入信号变得极其不稳定。因此，为了避免该情况的发生，需改进对输入信号 x 的规范方法，其计算为：

$$x^N = \frac{x - \bar{x}}{s} \quad (4)$$

式中， $s = \sqrt{\frac{\sum_{i=1}^N x_i^2 - N\bar{x}^2}{N-1}}$ 是标准偏差。

1.3 映射函数的求解

文献[5]中，映射函数是采用线性最小二乘法(LLS)对式(1)进行估计。该方法没有考虑噪声的影响，在对ECG的逆问题、神经网络等问题的研究中，也存在着类似的情况。在很多的相关研究中，已经证明采用奇异值截断(truncated singular value decomposition, TSVD)和总体最小二乘(total least squares, TLS)能有效地抑制噪声的干扰。由于LLS、TLS、STR和TSVD方法均为线性函数方法，它们所构建的线性映射函数不能很好地刻画误差与特征之间的关系，因此可以考虑采用非线性映射函数对DTS的预测。采用SVR方法构建非线性映射函数可以更好地刻画误差与特征之间的关系，本文研究

SVR方法所构建的映射函数对预测无损ICP的影响。

1.3.1 线性最小二乘法(LLS)

为了达到预测误差 e_i 的目的，在给定特征向量 F_i 的前提下，采用如下线性函数：

$$e_i = F_i^T b \quad (5)$$

式中， F_i 是血液动力学向量； e_i 代表输出的误差； b 是与血液动力学相关的向量，其维数和 F_i 一致，给定训练数据集得出的 e_i 和 F_i ，其中 $i = 1, 2, \dots, N_k$ ，对此，应用线性最小二乘法(LLS)可以得到估计值 b ：

$$\hat{b} = (F^T F)^{-1} F^T e \quad (6)$$

式中， F 是一个 $N_k \times d$ 的数据矩阵， e_i 是由第 i 行 $F_i^T e$ 构成的列向量构成， N_k 是输入的总数减去排除掉的一个数。根据标准线性回归理论可得出估计向量 \hat{b} 为：

$$V_b = (F^T F)^{-1} \sigma^2 \quad (7)$$

式中， σ^2 是 e 的噪声向量，假设其为高斯白噪声。

1.3.2 支持向量回归方法

支持向量回归方法(SVR)和许多基于风险最小化原则的常规方法相比较，主要差别是其使用结构风险最小化，追求模型预测性和泛化力，因此在解决小样本、非线性、高维数灾难和局部极小等问题时比常规方法有更好的表现。对于训练样本集 $\{x_i, y_i\}$ ， $x_i \in R^n$ 为输入变量的值， $y_i \in R$ 为相应的输出值。支持向量回归的基本思想是通过从输入空间到输出空间的非线性映射 ϕ ，将输入信号 x 映射到高维特征空间，并在特征空间中采用下述函数进行回归：

$$y = (\omega \phi(x)) + b \quad \phi: R^n \rightarrow F, \omega \in F \quad (8)$$

根据统计学习理论的结构风险最小化准则，支持向量回归方法通过极小化目标函数来确定如下回归函数式：

$$\min \frac{1}{2} w^T w + C \sum_{i=1}^l (\xi_i + \xi_i^*) \quad i = 1, 2, \dots, l \quad (9)$$

$$\text{s.t.} \begin{cases} y_i - w^T \phi(x_i) - b \leq \varepsilon + \xi_i \\ w^T \phi(x_i) + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases}$$

式中, 惩罚因子 C 为可调参数, C 越大对错误的惩罚越重。在考虑训练数据有噪音的情况时, 使用 C 控制经验风险。核函数中自由参数 l 的选择和惩罚因子 C 的选择会决定不同的模型(决策函数)。选择参数的值从而找到最小化风险的决策函数的过程称作模型选择。在实现的时候, 非线性映射采用如下的径向基高斯核函数:

$$k(x, x_i) = \exp\left(\frac{-\|x - x_i\|^2}{2\sigma^2}\right) \quad (10)$$

参数的选择采用试凑法与最小最大化相结合的原则:

- 1) 为常数和核函数固有参数赋初值;
- 2) 最大化得到 ω 和 b ;
- 3) 更新惩罚因子 C 和核参数, 最小化估计值;
- 4) 如果估计值满足要求, 结束运算; 否则重复2)。

1.4 查询策略

本文采用的查询策略为最小误差查询, 即选择有最小估计误差的模型作为理想的预测模型对ICP进行估计。

按照1.3节中的方式对线性映射进行估计以后, 从有 n 个训练样本的数据库中, 可获得 n 个不同的映射函数。对于一个输入向量, 基于这 n 个映射函数, 获得 n 个误差的估计值 $e_j, 1 \leq j \leq n$ 。

2 结果

在研究中所使用的数据采自某医学院的42个具有脑部疾病的病人, 每个病人数据记录长度为5 min~25 min, 每条记录长度大约100个心跳周期, 用作数据库的一个输入, 总共有242个输入数据对。ABP、ICP和CBFV等3种信号都从相关的采集设备用75 Hz的采样频率获得。由于只采用了42个病人的较小样本量, 因此在性能评估时, 基于Leaving-one-out 策略进行有效性的交叉验证: 首先是把所有样本分割成2部分, 在 N 个样本中取 $(N-1)$ 个样本作为训练样本; 其次把剩下的那个样本作为测试样本, 以验证有效性。重复 N 次该处理过程, 目的是让所有样本都能够被用作训练样本和测试样本, 能够得到更有效的利用。

采用基于LLS和SVR方法所估计的映射函数在估计误差分别为设定阈值的10%、25%、50%、75%及90%时的误差值如表1所示, 相应的ICP预测波形分别如图3和图4所示。

表1 对应于估计误差分别为设定阈值的10%、25%、50%、75%及90%时所得到的的误差值

方法	误差值/(%)				
	10	25	50	75	90
LLS	0.082	0.153	0.262	0.482	0.692
SVR	0.079	0.140	0.249	0.414	0.639

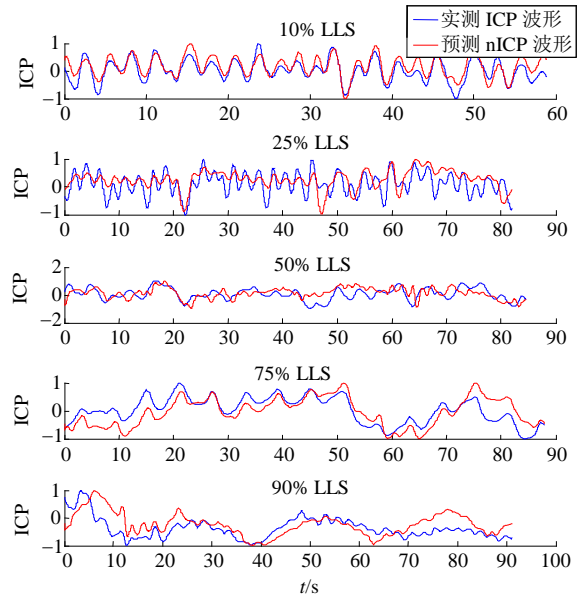


图3 基于LLS方法构成的映射函数的ICP波形

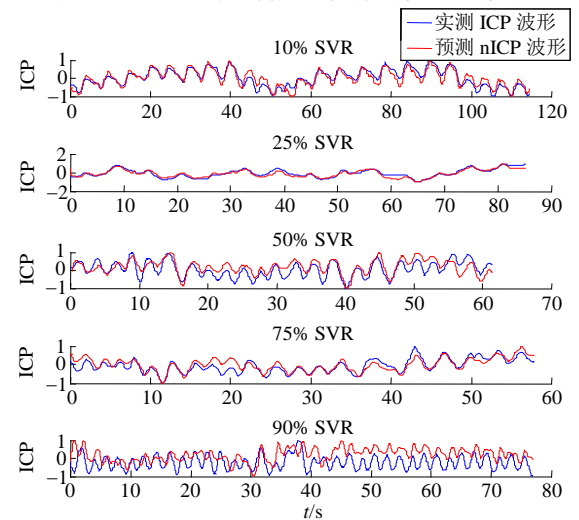


图4 基于SVR方法构成的映射函数的ICP波形

3 讨论和结论

表1显示不同映射函数构成的估计方法在实际应用中, 对无损ICP的预测具有较大的影响。当采用线性最小二乘法时, 90%的估计误差都小于0.692, 而采用支持向量回归方法(非线性)时, 90%的估计误差都小于0.639, 较前面的线性方法有较好的性能提升。图4中的结果直观地显示了基于这些映射函数对

ICP波形的预测结果,预测波形与真实波形有很好的 consistency,能够较好地捕获ICP信号中的瞬时变化情况,为用户提供直观的波形预测结果。这正是该时间系列数据挖掘方法相对于当前其他ICP无损估计方法的优点^[6-7]。

以上结果表明,当采用更为复杂的非线性方法对特征和误差间的关系进行描述,能够获得更为稳健的ICP估计结果。

参 考 文 献

- [1] CHEN C H, NEVO E, FETICS B, et al. Estimation of central aortic pressure waveform by mathematical transformation of radial tonometry pressure, validation of generalized transfer function[J]. *Circulation*, 1997, 95(7): 1827-1836.
- [2] FETICS B, NEVO E, CHEN C H, et al. Parametric model derivation of transfer function for noninvasive estimation of aortic pressure by radial tonometry[J]. *IEEE Trans Biomed Eng*, 1999, 46(6): 698-706.
- [3] SWAMY G, LING Q, LI T, et al. Blind identification of the aortic pressure waveform from multiple peripheral artery pressure waveforms[J]. *Am J Physiol Heart Circ Physiol*, 2007, 292(5): 2257-2264.
- [4] HU X, NENOV V, BERGSNEIDER M, et al. Estimation of

- hidden state variables of the intracranial system using constrained nonlinear kalman filters[J]. *IEEE Trans Biomed Eng*, 2007, 54(4): 597-610.
- [5] HU X, XU P, WU S Z, et al. A data mining framework for time series estimation[J]. *Journal of Biomedical Informatics*, 2010, 43(2): 190-199.
- [6] WU S, WU Y, XU P, et al. Compare time series mining approaches for mapping function assessment[C]/IN ICCAS2009. San Jose, USA: [s.n.], 2009.
- [7] WU S, XU P, ASGARI S, et al. Time series mining approach for noninvasive intracranial pressure assessment: an investigation of different regularization techniques[C]/CSIE. Los Angeles, USA: [s.n.], 2009: 382-386.
- [8] PANERAI R B, DAWSON S L, POTTER J F. Linear and nonlinear analysis of human dynamic cerebral autoregulation[J]. *American Journal of Physiology-Heart and Circulatory Physiology*, 1999, 277(3 Pt 2): 1089-1099.
- [9] OSUNA E, FREUND R, GIROSI F. Support vector machines: training and applications[R]. Cambridge, MA: MIT, 1997.
- [10] VAPNIK V, GOLOWICH S, SMOLA A. Support vector method for function approximation, regression estimation, and signal processing[M]. Cambridge, MA: MIT Press, 1997: 281-287.

编辑 蒋 晓

(上接第955页)

- [3] MEANEY P M, FANNING W M, RAYNOLDS T, et al. Initial clinical experience with microwave breast imaging in women with normal mammography[J]. *Academic Radiology*, 2007, 14(2): 207-218.
- [4] WINTERS D W, SHEA J D, HAGNESS S C, et al. Three-dimensional microwave breast imaging: dispersive dielectric properties estimation using patient-specific basis functions[J]. *IEEE Trans Med Imag*, 2009, 28(7): 969-981.
- [5] WADBRO E, BERGGREN M. High contrast microwave tomography using topology optimization techniques[J]. *Journal of Computational and Applied Mathematics*, 2009, doi:10.1016/j.cam.2009.08.027.
- [6] REIGBER A, LAURENT F F. Interference suppression in synthesized SAR images[J]. *IEEE Geoscience and Remote Sensing Letters*, 2005, 2(1): 45-49.
- [7] 董臻, 梁甸农, 黄晓涛. VHF/UHF UWB SAR基于通道均衡的RFI抑制方法[J]. *电子与信息学报*, 2008, 30(3): 550-553.

- DONG Zhen, LIANG Dian-nong, HUANG Xiao-tao. A RFI suppression algorithm based on channel equalization for the VHF/UHF UWB SAR[J]. *Journal of Electronics & Information Technology*, 2008, 30(3): 550-553.
- [8] 李廷军, 孔令讲, 周正欧. 随机相位编码抑制步进频率探地雷达RFI的研究[J]. *电子与信息学报*, 2009, 31(7): 1771-1774.
- LI Ting-jun, KONG Ling-jiang, ZHOU Zheng-ou. Researching on random phase-coded to suppress RFI in SFGPR[J]. *Journal of Electronics & Information Technology*, 2009, 31(7): 1771-1774.
- [9] CIOCAN R, JIANG H. Model-based microwave image reconstruction: simulations and experiments[J]. *Med Phys*, 2004, 31: 3231-3241.
- [10] JIANG Hua-bei, LI Chang-qing. Ultrasound-guided microwave imaging of breast cancer: Tissue phantom and pilot clinical experiments[J]. *Med Phys*, 2005, 32(8): 2528-2535.

编辑 黄 莘