

基于支持向量机回归的短时交通流预测模型^{*}

傅贵¹ 韩国强¹ 逯峰² 许子鑫¹

(1. 华南理工大学 计算机科学与工程学院, 广东 广州 510640; 2. 广州市交通管理科学技术研究所, 广东 广州 510640)

摘 要: 将交通流预测的理论和方法引入交通控制系统, 可提高交通控制系统对交通流变化的自适应能力. 为此, 文中通过引入核函数把短时交通流预测问题转化为高维空间中的线性回归问题, 提出了基于支持向量机回归的短时交通流预测模型, 并利用广州市交通流检测系统的数据进行实验. 结果表明, 文中模型的预测结果与实际数据相吻合, 预测误差小于基于卡尔曼滤波的预测方法, 从而验证了该模型的可行性和有效性.

关键词: 交通控制; 短时交通流; 预测模型; 机器学习; 支持向量机回归

中图分类号: TP391

doi: 10.3969/j.issn.1000-565X.2013.09.012

城市道路交通是一个动态的复杂系统, 随着观测时间范围的缩小, 交通特征由确定性向随机性过渡转化, 交通流的预测难度也随之提高. 近年来, 有关短时交通流预测理论的研究得到国内外众多学者的高度关注, 已成为当前智能交通的研究热点之一. 通常, 将未来 5 min 至 1 h 的交通流预测称为短时交通流预测. 短时交通流预测在智能交通控制和诱导方面具有广泛的应用前景, 可以缓解或解决城市交通拥堵问题. 短时交通流具有非线性、随机性和不确定性等特征. 目前, 城市交通控制和诱导普遍采用预设方案, 少数城市采用基于交通流检测的自适应控制模式, 基于短时交通流预测的智能交通控制和诱导的应用仍然较少.

近年来, 应用于短时交通预测的模型主要有统计模型、非线性预测模型、神经网络模型、组合模型和微观交通仿真模型等. 对于统计模型, Ahmed 等^[1]首次将时间序列模型应用于交通流预测领域; Vythoulkas^[2]提出了基于卡尔曼滤波的交通流预测

模型, 并应用于驾驶员信息服务系统中. 对于非线性预测模型, Dendrinis^[3]将突变理论应用于交通数据分析中; Huang 等^[4]基于相空间重构理论提出了城市交通流的非线性混沌预测模型. 对于神经网络模型, Corinne^[5]利用神经网络建立了每个路段及整个路网的交通流预测模型, 并利用模拟数据对模型进行验证, 取得了较好的预测效果. Jiang 等^[6]提出了一种用以预测实际交通流的动态小波神经网络模型. 邵春福等^[7]提出了基于支持向量机回归的交通状态短时预测方法.

为提高交通控制系统对交通流变化的自适应能力, 文中提出了基于支持向量机回归 (SVMR) 的短时交通流预测模型, 并采用广州市交通流检测系统的数据对模型进行实验和定量分析, 以验证模型的可行性、有效性.

1 支持向量机

支持向量机 (SVM) 通过结构风险最小化较好

收稿日期: 2013-02-05

^{*} 基金项目: NSFC-广东省政府联合基金资助项目 (U1035004); 国家自然科学基金青年科学基金资助项目 (61003270); 广州市科技计划重点支撑项目 (11A11080267); 广东省计算科学重点实验室开放基金资助项目 (201206005)

作者简介: 傅贵 (1975-), 男, 在职博士生, 高级工程师, 主要从事智能交通系统技术研究. E-mail: longman@188.com

地解决了小样本、非线性、维数灾难、过学习和局部极小等问题,已成为机器学习领域的研究热点之一^[8]. SVM 的主要思想是:给定训练样本,建立一个超平面作为决策曲面,使得正例和反例之间的间隔最大化^[9]. SVM 可用于求解模式识别和非线性回归问题.

1.1 最优超平面

超平面是从线性可分模式的线性分类器发展而来的. 在二维空间中,如果一个线性函数能将样本数据完全分开,则称样本数据是线性可分的,该分类线性函数可表示为

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad (1)$$

式中, \mathbf{x} 为样本向量, \mathbf{w} 为样本向量的法向量, b 为偏移常量.

在图 1 中, H_1 、 H_2 平行于 H 且可区分各类样本. H_1 、 H_2 的训练样本点称为支持向量(SV), H_1 、 H_2 之间的距离称为分类间隔. 最优超平面就是使分类间隔最大的平面.

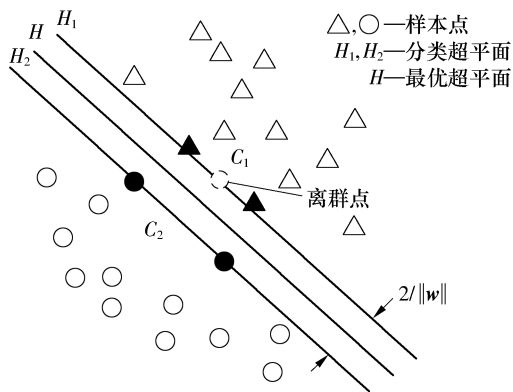


图 1 最优超平面示意图

Fig. 1 Schematic diagram of optimal hyperplane

第 i ($i = 1, 2, \dots, N$; N 为训练样本数) 个训练样本由一个向量和一个类别组成,可表示为

$$D_i = (\mathbf{x}_i, y_i) \quad (2)$$

式中, \mathbf{x}_i 为输入向量, y_i 为类别标记. 对于二元线性分类, y_i 只有 1 和 -1 两个值. 样本点到某个超平面的间隔为

$$\delta_i = y_i(\mathbf{w}^T \mathbf{x}_i + b) = |g(\mathbf{x}_i)| \quad (3)$$

将式(1)中 \mathbf{w} 和 b 进行归一化处理,可得到点到超平面的欧氏距离(即几何间隔):

$$\delta_i = \frac{1}{\|\mathbf{w}\|} |g(\mathbf{x}_i)| \quad (4)$$

将所有样本点中最小的间隔定义为 1, 此时相

应的两条极端直线的几何间隔为 $\frac{2}{\|\mathbf{w}\|}$, 如图 1 所示.

由于几何间隔与 $\|\mathbf{w}\|$ 成反比, 最大化分类间隔可以通过最小化 $\|\mathbf{w}\|$ 获得, 即

$$\min \|\mathbf{w}\| \quad (5)$$

等价于

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad (6)$$

文中规定样本点必须在 H_1 或 H_2 的一侧或者在 H_1 、 H_2 上, 由于所有样本点之间的间隔大于 1, 因此

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0, \quad i = 1, 2, \dots, N \quad (7)$$

最大分类间隔的求解等价于在约束条件 $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0$ 下使 $\frac{1}{2} \|\mathbf{w}\|^2$ 最小.

1.2 松弛变量

在实际问题中, 训练集也可能会出现不可分的样本点(称为离群点, 如图 1 所示), 它们会影响分类超平面的形成^[10].

为处理不可分离的数据点, 文中引入松弛变量 $\{\xi_i\}_{i=1}^N$, 用于度量数据点对模式可分理想条件的偏离程度. 当 $0 < \xi_i \leq 1$ 时, 数据点落入超平面的正确一侧; 当 $\xi_i > 1$ 时, 数据点落在超平面的错误一侧. 因此, 不可分问题可以描述为: 给定训练样本 $\{(\mathbf{x}_i, d_i)\}_{i=1}^N$, 寻找权值向量 \mathbf{w} 和偏置 b , 使得

$$\Phi(\mathbf{w}, \boldsymbol{\xi}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \quad (8)$$

极小, 且满足约束条件

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad (9)$$

$$(i = 1, 2, \dots, N; \xi_i \geq 0)$$

其中, C 为用户选定的正参数, $\boldsymbol{\xi} = (\xi_1, \xi_2, \dots, \xi_N)$.

2 建模过程

SVMR 是支持向量机在回归估计问题中的扩展, 解决支持向量机回归问题的目标是: 让所有样本点逼近超平面, 使得样本点与超平面的总偏差达到最小. 应用于交通流预测的支持向量机回归主要有 ε -支持向量机回归、 v -支持向量机回归和最小二乘支持向量机回归(LS-SVMR)等^[10]. 文中拟将支持向量机回归应用于交通流预测中. 短时交通流预测属于非线性回归问题, 文中通过引入核函数, 把短时交通流预测问题转化为高维空间(Hilbert 空间)中的线性回归问题^[11].

2.1 ε -支持向量机回归

如果空间 \mathbf{R}^n 中存在超平面 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, 使得 $|f(\mathbf{x}) - y| \leq \varepsilon (\forall (\mathbf{x}_i, y_i) \in S, \varepsilon > 0)$, 则称 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ 是样本集合 S 的 ε -线性回归^[12].

支持向量机回归要解决的问题为

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad (10)$$

$$\text{s. t. } |\mathbf{w}^T \mathbf{x}_i + b - y_i| \leq \varepsilon, \quad i = 1, 2, \dots, N.$$

引入松弛变量, 式(10)可写为

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*) \quad (11)$$

$$\text{s. t. } \begin{cases} (\mathbf{w}^T \mathbf{x}_i + b) - y_i \leq \varepsilon + \xi_i \\ y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0 \\ \xi_i^* \geq 0 \\ i = 1, 2, \dots, N \end{cases}.$$

式中, ξ_i^* 为松弛变量.

对最优化问题(11)采用拉格朗日乘子法, 得到

$$\min \left\{ L = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*) - \sum_{i=1}^N \alpha_i (\varepsilon + \xi_i - y_i + \mathbf{w}^T \mathbf{x}_i + b) - \sum_{i=1}^N \alpha_i^* (\varepsilon + \xi_i^* + y_i - \mathbf{w}^T \mathbf{x}_i - b) - \sum_{i=1}^N (\beta_i \xi_i + \beta_i^* \xi_i^*) \right\} \quad (12)$$

式中, $\alpha_i, \alpha_i^*, \beta_i, \beta_i^*$ 为拉格朗日乘子.

求解式(12)并进行对偶变换, 可得到

$$\min \left\{ \frac{1}{2} \sum_{i,j=1}^N (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^N (\alpha_i^* - \alpha_i) y_i + \sum_{i=1}^N (\alpha_i + \alpha_i^*) \varepsilon \right\} \quad (13)$$

$$\text{s. t. } \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0, \quad 0 \leq \alpha_i, \alpha_i^* \leq C.$$

2.2 核函数的选择

对于非线性问题, 超平面已经无法进行分类. 如图2所示, 若将点 m 和 m' 之间所有的点定义为正

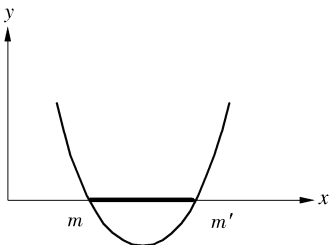


图2 超平面对非线性问题的错误分类示例

Fig.2 Example of misclassification about hyperplane in nonlinear problem

类, 其余两边的点定为负类, 则无法用一个线性函数把两类正确地分开, 但可通过点在曲线的上方还是下方来判断点的类别.

对于线性不可分的 $g(\mathbf{x})$, 文中通过核函数 $K(\mathbf{x}, \mathbf{x}^*)$ 将线性不可分问题转化成高维空间的线性问题, $K(\mathbf{x}, \mathbf{x}^*)$ 接收低维空间的输入值, 能计算出高维空间的内积值. SVM 的结构类似于一个神经网络^[9], 中间层的每个节点 $K(\mathbf{x}_i, \mathbf{x})$ 对应于一个支持向量, 通过核函数的变化得到中间层节点的线性组合的输出结果, 如图3所示.

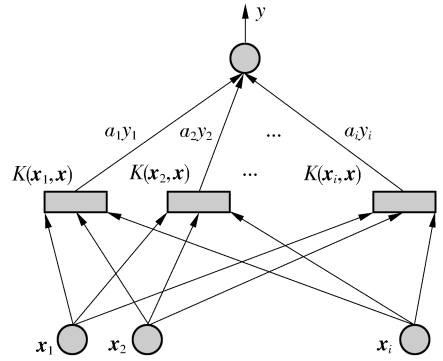


图3 使用核函数的支持向量机结构

Fig.3 Structure of SVM with kernel function

对于非线性回归问题, 文中将内积用核函数替代, 对式(13)引入核函数, 可得

$$\min \left\{ \frac{1}{2} \sum_{i,j=1}^N (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N (\alpha_i^* - \alpha_i) y_i + \sum_{i=1}^N (\alpha_i + \alpha_i^*) \varepsilon \right\} \quad (14)$$

$$\text{s. t. } \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0, \quad 0 \leq \alpha_i, \alpha_i^* \leq C.$$

求解二次规划问题(14)可得到 α_i 和 α_i^* , 并由

此得到 $\mathbf{w} = \sum_{i=1}^N (\alpha_i^* - \alpha_i) \mathbf{x}_i$ 和 b , 即

$$b = \begin{cases} y_j + \varepsilon - \sum_{i,j=1}^N (\alpha_i^* - \alpha_i) K(\mathbf{x}_i, \mathbf{x}_j), & 0 \leq \alpha_i \leq C \\ y_j - \varepsilon - \sum_{i,j=1}^N (\alpha_i^* - \alpha_i) K(\mathbf{x}_i, \mathbf{x}_j), & 0 \leq \alpha_i^* \leq C \end{cases} \quad (15)$$

选择适当的核函数是关键. SVM 普遍使用3类核函数^[13]: 多项式核函数 $K(\mathbf{x}_i, \mathbf{x}) = [\langle \mathbf{x}_i, \mathbf{x} \rangle + 1]^q$ 、高斯径向基核(RBF)函数 $K(\mathbf{x}_i, \mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}\|^2}{\sigma^2}\right)$ 和两层感知器核函数 $K(\mathbf{x}_i, \mathbf{x}) = \tanh(v \langle \mathbf{x}_i, \mathbf{x} \rangle + c)$, 文中使用 RBF 函数.

2.3 建模流程

假设 $\mathbf{x}_i (\mathbf{x}_i \in \mathbf{R}^n)$ 为影响交通流预测的因素, y_i 为交通流预测值. 基于 SVMR 的短时交通流预测模型就是寻找 \mathbf{x}_i 与 y_i 之间的关系^[14]:

$$f: \mathbf{R}^n \rightarrow \mathbf{R} \quad (16)$$

$$y_i = f(\mathbf{x}_i) \quad (17)$$

文中采用当前 t 和前 n 个时段的交通流作为输入值, 对未来 $t+1$ 时段的交通流进行预测^[15]. 具体建模步骤如下:

(1) 选择样本数据并做去噪归一化等预处理, 构造训练集.

设当前时段的交通流为 $q_i(t)$, 则下一时段 $q_i(t+1)$ 的训练样本集为 $\mathbf{x}_i = (q_i(t), q_i(t-1), \dots, q_i(t-n))$.

(2) 通过对已知数据的分析选择核函数 $K(\mathbf{x}_i, \mathbf{x}) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}\|^2}{\sigma^2}\right)$ 和适当的参数.

(3) 利用训练样本建立目标函数, 通过求解二次规划问题(14)来寻找最优分类面.

(4) 由求得的最优解构建决策函数, 即

$$f(\mathbf{x}) = \sum_{i=1}^N (\alpha_i^* - \alpha_i) K(\mathbf{x}_i, \mathbf{x}) + b \quad (18)$$

用测试样本集计算未来时刻的预测值.

图4为基于 SVMR 的短时交通流的建模流程^[10].

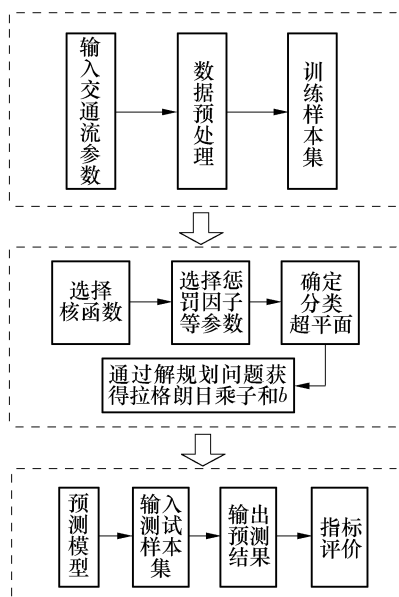


图4 基于 SVMR 的短时交通流的建模流程图^[10]

Fig.4 Flowchart of short-term traffic flow modeling based on SVMR

3 实例验证

3.1 数据来源和计算结果

为验证所提预测模型的有效性, 文中使用广州市交通流检测系统的真实数据进行实验. 对东风西路南往北(断面1)、解放路南往北(断面2)和科韵路南往北(断面3)3个断面进行采样, 采样间隔为5 min, 将2011-10-27T05-55-00—23-50-00的212条流量数据作为基于 SVMR 的预测模型的训练样本集, 2011-10-28T05-55-00—19-55-00的164条流量数据作为测试样本集.

采用文中提出的模型对断面1、2、3的交通流进行预测. 其中, SVM 核函数选用 Gauss RBF 函数, 其参数 σ 的经验公式为 $\sigma = k^{-1}$, k 为输入数据中属性的个数. 文中将当前时段的流量和前4个时段的流量作为输入值, 属性总个数为5, 故 $\sigma = 0.2$. 对于惩罚因子 c , 较大的惩罚因子可以提高预测结果的准确率, 但过大的惩罚因子会造成过学习状态, 从而影响最终测试样本集的准确率. 根据经验设 $c = 0.8$, 损失函数 $\varepsilon = 0.1$. 采用训练集分别训练各断面的预测模型, 然后分别预测各断面在2011-10-28的164个时段的交通流量, 预测结果如图5所示, 文中模型的预测结果与实际数据吻合良好.

模型预测的交通流绝对误差如图6所示, 绝对误差量为真实数据 Y_i 与预测数据 Y_i^* 之差. 从图可知, 基于 SVMR 的预测模型能够有效地预测短时交通流, 每一断面的预测值的浮动范围不大, 预测效果良好.

3.2 效果分析

为定量分析文中模型的预测效果, 文中引入2个评价指标, 即平均平方误差(MSE)和平均绝对百分比误差(MAPE), 并将文中方法与卡尔曼滤波法的预测误差进行了比较. 卡尔曼滤波法是一种线性回归的预测方法, 目前已经成功应用于交通需求预测以及短时交通流预测. MSE 定义为

$$MSE = \frac{1}{N} \sum_{i=1}^N (Y_i - Y_i^*)^2 \quad (19)$$

MAPE 定义为

$$MAPE = \frac{1}{N} \sum_{i=1}^N \left| \frac{Y_i - Y_i^*}{Y_i} \right| \times 100\% \quad (20)$$

两种方法对3个断面交通流的预测误差如表1所示.

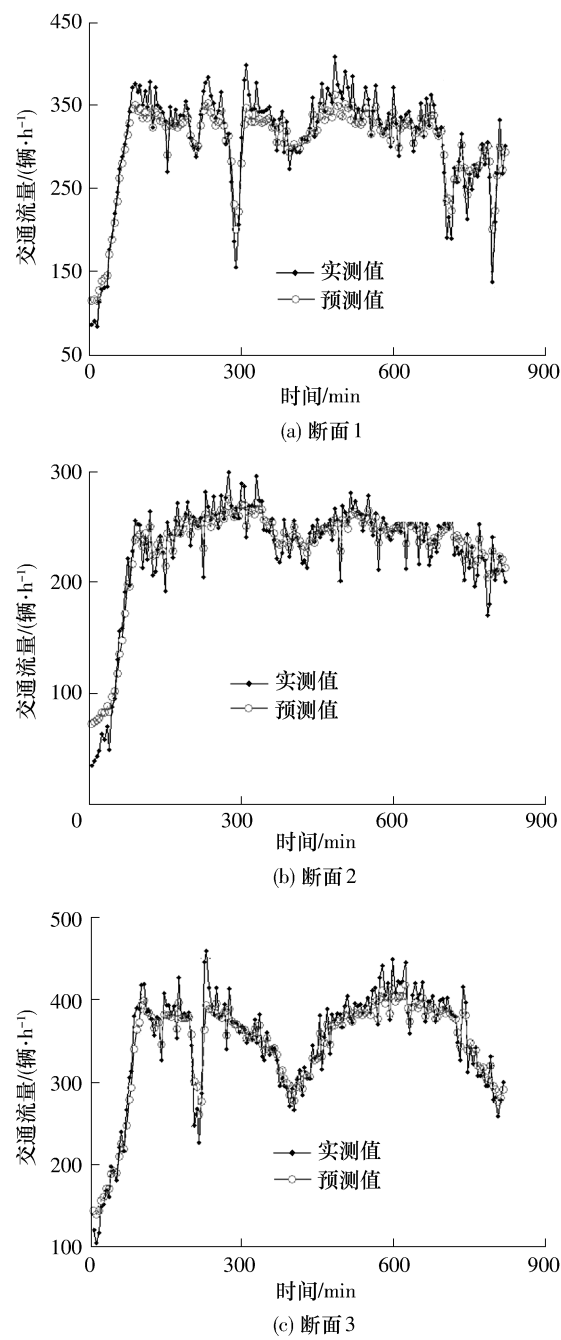


图5 3个断面的交通流量预测值与实测值对比

Fig.5 Comparison between forecasting values and measured values of traffic flow in 3 sections

表1 文中方法与卡尔曼滤波法的预测误差比较

Table 1 Comparison of forecasting errors between the proposed method and Kalmen filtering method

预测方法	断面	MAPE/%	MSE
卡尔曼滤波法	断面1	6.04	416.32
	断面2	9.50	212.91
	断面3	5.17	313.43
文中方法	断面1	5.57	362.43
	断面2	6.85	166.03
	断面3	3.81	284.88

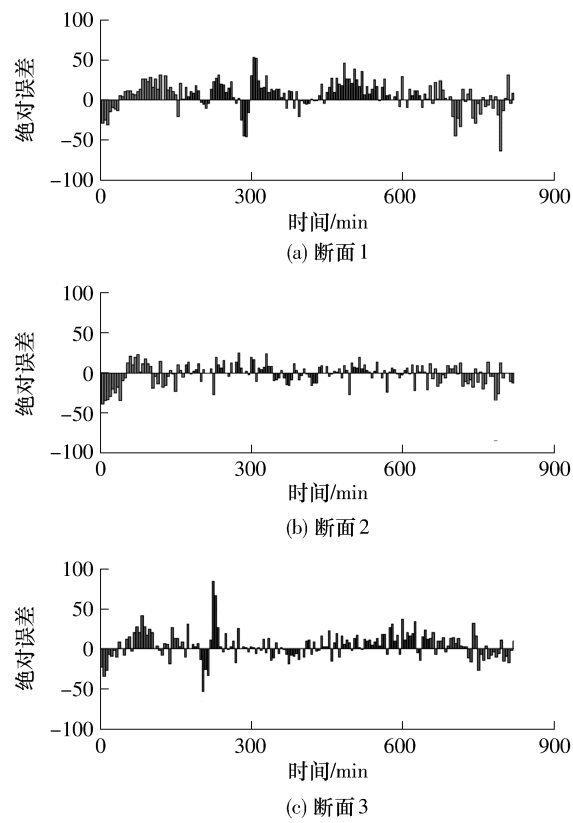


图6 3个断面的交通流量绝对误差

Fig.6 Absolute errors of traffic flow in three sections

表1表明,文中预测方法的 MAPE 和 MSE 均优于卡尔曼滤波法。

4 结语

文中建立了基于 SVMR 的短时交通流预测模型,其中选择合适的核函数及其参数(包括惩罚因子)是建模过程的关键.实证分析表明,文中建立的模型是可行和有效的.与卡尔曼滤波预测方法相比,文中模型的预测性能更优、更实用.适当的参数优化将有效地提高基于 SVMR 的预测准确性,今后将在参数优化方面进行深入研究。

参考文献:

[1] Ahmed M S, Cook A R. Analysis of freeway traffic time-series data by using box-Jenkins techniques [J]. Transportation Research Record, 1979, 722: 1-9.

[2] Vythoulkas P C. Alternative approaches to short term traffic forecasting for use in driver information systems [M]. Berkeley: Elsevier Science Publishers, 1993.

[3] Dendrinios D S. Operating speeds and volume to capacity

- rations; the observed relationship and the fold catastrophe [J]. *Transportation Research*, 1978, 12(3): 191-194.
- [4] Huang K, Chen S, Zhou Z G. Research on non-linear chaotic prediction mode for urban traffic flow [J]. *Journal of Southeast University: English Edition*, 2003, 19(4): 411-413.
- [5] Ledoux Corinne. An urban traffic flow model integrating neural network [J]. *Transportation Research Part C: Emerging Technologies*, 1997, 5(5): 287-300.
- [6] Jiang X, Adeli H. Dynamic wavelet neural network model for traffic flow forecasting [J]. *Journal of Transportation Engineering*, 2005, 131(10): 771-779.
- [7] 姚智胜, 邵春福, 高永亮. 基于支持向量回归机的交通状态短时预测方法研究 [J]. *北京交通大学学报*, 2006, 30(3): 19-22.
- Yao Zhi-sheng, Shao Chun-fu, Gao Yong-liang. Research on methods of short-term traffic forecasting based on support vector regression [J]. *Journal of Beijing Jiaotong University*, 2006, 30(3): 19-22.
- [8] 邵春福, 熊志华, 姚智胜. 道路网短时交通需求预测理论、方法和应用 [M]. 北京: 清华大学出版社, 2011: 58-61, 129-146.
- [9] Haykin Simon. 神经网络与机器学习 [M]. 3 版. 申富饶, 徐烨, 郑俊, 译. 北京: 机械工业出版社, 2011: 144-192.
- [10] 王凡. 基于支持向量机的交通流预测方法研究 [D]. 大连: 大连理工大学计算机学院, 2010.
- [11] Courant R, Hilbert D. *Methods of mathematical physics* [J]. *Physics Today*, 1954, 7(5): 17.
- [12] Vapnik V. 统计学习理论的本质 [M]. 张学工, 译. 北京: 清华大学出版社, 2000.
- [13] 张学工. 关于统计学习理论与支持向量机 [J]. *自动化学报*, 2000, 26(1): 32-42.
- Zhang Xue-gong. Introduction to statistical learning theory and support vector machines [J]. *Acta Automatica Sinica*, 2000, 26(1): 32-42.
- [14] Lin Chih-Jen. LIBSVM: a library for support vector machines [CP/OL]. (2013-04-01) [2013-04-20]. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>.
- [15] Wu C H, Wei C, Chang M. Travel time prediction with support vector regression [J]. *IEEE Transactions on Intelligent Transportation Systems*, 2004, 5(12): 276-281.
- [16] 杨兆升, 王媛, 管青. 基于支持向量机方法的短时交通流量预测方法 [J]. *吉林大学学报: 工学版*, 2006, 36(6): 881-884.
- Yang Zhao-sheng, Wang Yuan, Guan Qing. Short-term traffic flow prediction method based on SVM [J]. *Journal of Jilin University: Engineering and Technology Edition*, 2006, 36(6): 881-884.

Short-Term Traffic Flow Forecasting Model Based on Support Vector Machine Regression

Fu Gui¹ Han Guo-qiang¹ Lu Feng² Xu Zi-xin¹

(1. School of Computer Science and Engineering, South China University of Technology, Guangzhou 510640, Guangdong, China;

2. Guangzhou Research Institute of Traffic Management Science, Guangzhou 510640, Guangdong, China)

Abstract: As the short-term traffic flow forecasting theories and approaches help to improve the ability of traffic control systems to automatically adapt to traffic flow changes, this paper proposes a short-term traffic flow forecasting model based on the support vector machine regression by using a kernel function to transform the issues into a linear regression problem in Hilbert Space. Then, the corresponding experiments are conducted based on the data from the traffic flow detection systems in Guangzhou. It is found that the forecasted results accord well with the actual data, and that the forecasting error of the proposed model is less than those of the prediction methods based on Kalman filtering. Thus, the feasibility and effectiveness of the proposed model are verified.

Key words: traffic control; short-term traffic flow; forecasting model; machine learning; support vector machine regression