

一类基于 SVM/RBF 的气象模型预测系统

罗泽举¹, 宋丽红², 薛宇峰², 朱思铭¹

(1. 中山大学数学与计算科学学院, 广州 510275; 2. 湛江海洋大学海滨校区, 湛江 524005)

摘要: 利用广东湛江地区近 30 年月平均气温的气象数据作为数据集, 建立了一种新的基于径向基核函数的支持向量机模型预测系统。通过适当选择模型参数, 其平均绝对百分比误差只有 5.61%, 在绝对误差温度小于等于 2 的条件下, 预测的准确率达到 95%, 显示出所建立的支持向量机模型预测系统的有效性。通过分析发现了湛江海岸地区气候和全球气候变暖一致的事实。

关键词: 径向基核函数; 支持向量机; 模型预测系统

One Kind of Weather Model Forecast System Based on SVM/RBF

LUO Zeju¹, SONG Lihong², XUE Yufeng², ZHU Siming¹

(1. School of Mathematics and Computer Science, Sun Yat-Sen University, Guangzhou 510275;

2. Seashore Campus of Zhanjiang Ocean University, Zhanjiang 524005)

【Abstract】 Using the weather data of average temperature of nearly thirty years recorded at Guangdong Zhanjiang seashore region as the data set, this paper sets up a new support vector machines models forecast system based on radial basis kernel function(RBF). By choosing the model parameter properly, the mean absolute percentage error is only about 5.61%, if the absolute error is less than 2, the rate of accuracy comes to 95%. This shows the validity of the SVMS forecast system. In addition, by analysing it finds that the climate gets warm gradually in Zhanjiang and this is consistent with the global climate.

【Key words】 Radial basis kernel function; Support vector machines; Model forecast system

湛江位于我国大陆最南端、广东省西南部, 东濒南部海域, 西临北部湾, 背靠大西南。市区位于雷州半岛东北部, 东经 110°24', 北纬 21°12', 地处北回归线以南的低纬度地区, 属于亚热带海洋性季风气候, 冬无严寒, 夏无酷暑, 偶有台风天气, 四季宜人。近年来, 全球气候不断变暖, 湛江地区的气候不但发生着类似的变化, 而且又显示出湛江所独有的地方性特点。如何准确地分析和预测各种气象因素, 避免恶劣天气的影响, 对湛江地区经济发展具有重要意义。

本文利用湛江市 1951 年~2000 年的月平均气温资料, 分析了近 50 年来湛江地区温度的变化特点, 并利用支持向量机回归模型, 建立了温度变化预报模型, 对温度进行了年以上的预测, 并得到良好的预测效果。

1 支持向量机回归模型

1.1 回归支持向量机

利用支持向量机回归学习的函数是:

$$f(x) = w \cdot x + b \quad (1)$$

通过 Lagrange 乘子办法^[1], 估计出回归参数为

$$\bar{w} = \sum_{i=1}^l (\alpha_i - \alpha_i^*) x_i; \quad \bar{b} = -\langle \bar{\alpha}, (x_i + x_j) \rangle / 2 \quad (2)$$

式中 x_i, x_j 分别为两类支持向量。

若采用内积回旋的 SVM 机, 则相应的回归函数变形为

$$f(x, v, \beta) = \sum_{i=1}^l \beta_i K(x, v_i) + b \quad (3)$$

其中 $\beta_i, i=1, \dots, l$ 是标量, $v_i, i=1, \dots, l$ 是向量, $K(u, v)$ 为满足 Mercer 条件的函数^[2,3] (通常称为核函数), 本文选取的核函数为径向基核函数。形式为

$$K(x, x_i) = \exp\{-\gamma \|x - x_i\|^2\} \quad (4)$$

这是一个非常重要的参数项, 只有依赖于样本数据结构的核函数, 才能达到最佳的预测效果。

1.2 模型中几个重要参数的分析

(1)核函数 $K(u, v)$ 。根据 Mercer 定理, 只要是满足相关条件的实对称函数都可作为核函数。由于气象数据点呈非线性相关, 因此可以选取多种函数作为核函数。通过实验, 发现径向基核函数作为核函数效果最佳。最佳效果体现了核函数和样本数据本身的固有规律。

(2)损失函数和敏感度。选取二次损失函数的计算时间较少, 但是这并不表示用二次损失函数计算精确度高, 相反, 选取不敏感函数和合适的 γ 值可以达到最理想的预测值, 因为二次损失函数是不可调整精度的。从这点来讲, 不敏感函数比传统的二次损失函数优越。

(3)控制上界 C 。影响支持向量数目和计算时间。 C 值越小支持向量数目和计算时间越少, C 值不能太小, 否则误差会增大。通过训练, 得到以下模型参数:

表 1 模型参数的选取

参数类别	选取
核函数 $K(u, v)$	径向基函数
损失函数 L	不敏感函数
控制上界 C	$C=10$
敏感度	0.81

基金项目: 国家自然科学基金资助项目 (10371135)

作者简介: 罗泽举(1965 -), 男, 博士生, 主研方向: 机器学习与模式识别, 生物信息学; 宋丽红, 实验师; 薛宇峰, 讲师; 朱思铭, 教授、博导

收稿日期: 2005-11-28 **E-mail:** luozeju@126.com

2 温度序列数据分析

图1是湛江地区近50年以来历年的平均气温、5年滑动平均以及年平均温度的距平曲线分布图。由图可见,地处亚热带区域的湛江地区近50年来年平均气温基本比较稳定,50年平均温度为23.22;除1973年外,1967年~1985年的近20年间,温度距平都低于0,属于一段气候偏冷期(1976年达到最低温度22.45)。但1986年以后,基本上都成为正距平,平均温度呈准线性上升。1998年到达最高24.48,比多年平均值高出2.03,成为异常偏暖年份。

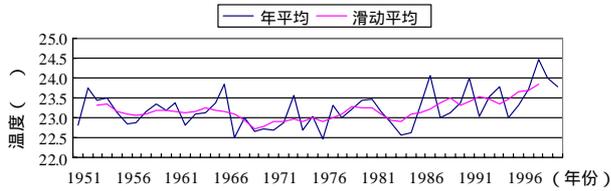


图1 湛江地区1951年~2000年历年平均气温、滑动平均曲线

分析1951年~2000年多年月平均气温可知,湛江地区各月平均温度变化幅度不是很大,因此全年温差不大,最冷月为1月(平均温度为15.8),最热月份为7月(平均温度达到28.9)。分析1月和7月平均气温距平,7月平均温度变化比较平缓,而1月气温变化的幅度就要大一些,而且自90年代以来,1月份气温的正距平持续偏大,说明冬季温度的变率比夏季气温的变率要大,而且近10来年冬季的增温也较夏季明显。通过计算1951年~1985年间的35年和1986年~2000年间的15年冬季(以12月、1月、2月平均温度作代表)和夏季(以6月、7月、8月平均温度为代表)的季节平均温度,可以看到,湛江地区冬季增温(0.8)较夏季(0.5)更明显。为此,我们分别分析了1951年~1985年和1986年~2000年两段时期内的月平均温度(见表2)。可见,月平均温度后15年比前35年有所升高,平均升高约0.65,这和近年来全球气候变暖的事实相一致。

表2 冬、夏季节平均温度反映的气温变化比较

季节	冬季	夏季
1951~1985	16.3	28.4
1986~2000	17.1	28.9

图2为1951年~1985和1986年~2000两个时间段内月平均气温对比图,红线代表1951年~1985年月平均温度,蓝线代表1986年~2000年各月平均气温,从中可以看出80年代中期以后气候变暖。

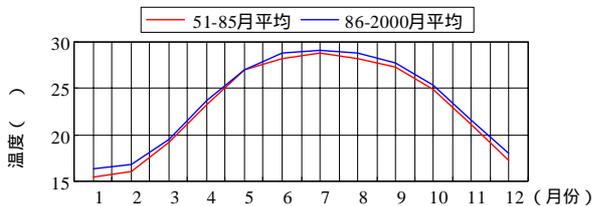


图2 月平均气温对比

另外,分析年际温度的变化发现,20世纪80年代之前,湛江地区的平均温度变化幅度不大,甚至50~70年代的10年平均温度还略有下降。但80年代之后显著增温,80年代比70年代的气温增加了0.13/10年,90年代又比80年代增加了0.51/10年,远远高于相应时段内全国平均温度增长水平^[4],这是地处亚热带海洋性季风气候的湛江地区平均温度的显著特点。

3 决策函数的确定

温度是一个时间序列 $x(t)$, $t=1, 2, \dots$, 其 t 时刻的值 $x(t)$

可以表示为函数:

$$x(t) = F(x(t-1), x(t-2), \dots, x(t-n)) \quad (5)$$

这里关键是要确定窗口参数 n 的值,以30年的训练数据为基础,假设当前值和前面3~7个时延相关(虽然理论上时延还可以有更宽的选择),限定窗口范围在3~7,计算得到平均绝对误差(MAE)如表3。

表3 模型输入窗口参数的确定

参数 n	3	4	5	6	7
MAE	1.7	1.89	5.69	4.30	3.21

其中最小的平均绝对误差是1.7,因此选取时延窗口(相空间维数)大小 $n=3$,即确定输入模型为

$$x(t) = F(x(t-1), x(t-2), x(t-3)) \quad (6)$$

4 预测结果分析

我们以时延为3进行预测,也就是用过去的3个数据预测将来的值,其平均绝对误差(MAE)、平均绝对百分比误差(MAPE)、根方差(RMSE)及相关系数(r)如表4。

表4 4个重要的预测指数

MAE	MAPE	RMSE	r
1.2163	5.61%	1.5358	0.98

可见其平均相对误差只有5.61%,根方差也相当小,只有1.5358,而1991年~2000年的10年间各年数据呈显著相关性,达0.98,说明所建立的预测模型是相当有效的。

图3是模拟曲线图和预测曲线。

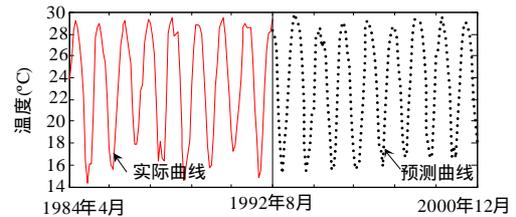


图3 1984年4月~2000年12月各月模拟曲线和预测曲线

由图可见,1984年4月~1992年8月各月实际和预测部分模拟相当吻合。以1992年9月~2000年12月作为预测,1984年4月~2000年12月各月误差曲线如图4所示。

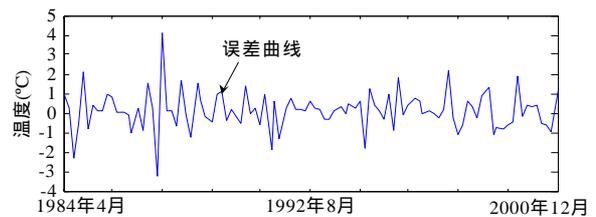


图4 1984年4月~2000年12月各月误差曲线

我们用3种核函数进行模拟预测,发现当选取径向基核函数时,其中误差大于2只占5%,最大误差为4,但只占3%。其它两种核函数的准确性方面都不如径向基核函数(见表5)。

表5 支持向量机模型:1992年9月~2000年12月各月预测准确率

误差温度限	预测月数(个)	预测准确率		
		RBF核	线性核	样条核
1	100	78%	70%	68%
2	100	95%	90%	91%
3	100	97%	92%	92%

当我们试用神经网络进行相同的测试时,选用径向基神经网络、线性神经网络、BP神经网络,其总体预测水平平均不如支持向量机模型(表6)。这表明支持向量机处理小样本的能力比传统神经网络预测能力要强。(下转第93页)