文章编号:1000-8608(2012)01-0079-06

基于相关向量机的中长期径流预报模型研究

仕玉治^{1,2},彭勇^{*1},周惠成¹

(1.大连理工大学水利工程学院,辽宁大连 116024;2.山东省水利科学研究院,山东济南 250013)

摘要:鉴于其优越的预报性能,将相关向量机(RVM)应用到中长期径流预报中,并在相空间重构的基础上,建立了基于相关向量机的径流预报模型.该模型首先对径流时间序列进行 相空间重构,并以重构后的径流序列作为模型输入;其次,采用粒子群优化(PSO)算法识别模型参数,利用优化所得重构参数验证时间序列具有混沌特性,在模型内循环过程中采用 EM 算法迭代估计超参数,并将 RVM 与应用较为广泛的最小二乘支持向量机(LSSVM)和自动 回归滑动平均模型(ARMA)进行了比较分析,结果表明该模型具有较好的泛化能力;最后, 基于水文过程变化的不确定性、RVM 描述输出值的不确定度以及相应概率下的预报区间, 使得调度人员在决策中能考虑预报的不确定性,定量估计各种决策的风险和效益.

关键词:相空间重构;相关向量机;长期径流预报;PSO 算法 **中图分类号:** TV124 **文献标志码:** A

0 引 言

中长期径流预报对统筹安排防洪与抗旱、水库 调度与管理等事务,实现水资源的最大利用效益具 有十分重要的实践意义,由于水文系统本身的非线 性及水文要素变化的不确定性,目前基于严密的物 理方法还很难对径流等水文现象进行描述和预测, 人们主要借助于成因分析法、水文统计法、模糊分 析等方法来描述和预测水文过程,水文统计法依据 水文资料的统计规律进行预测,方法较为常用,它 包括两大类:一类是分析水文要素自身随时间变化 的统计规律,并建立模型进行预测,如时间序列分 析法:另一类是分析水文要素与多因子之间的相关 关系,建立模型进行预测,如多元回归法.两类均可 直接利用原序列建立线性或非线性关系进行预测, 但其精度有时不能满足工程需要.相空间重构成为 分析时间序列的一种崭新的方法,通过挖掘或者恢 复水文系统的多变量影响因子,重构水文非线性动 力系统,国外许多学者对短期径流预报进行研究, 并取得了较好的应用效果^[1,2].时间滞时τ和嵌入 维数 m(文中也称重构参数)对时间序列的噪声和 数据量大小等影响因素比较敏感[3],通常采用互信 息法、关联维数法、虚假近邻法和 Cao 法等多种方法所得到的估计值差别较大,不利于获得较好的预 报精度.本文采用 Yu 等^[4]和Sivakumar^[5]提出的方法,优化得到重构参数.

一般线性方法难以描述水文系统非线性特征, 许多新的方法逐步被引用到水文预报模型当中,如 贝叶斯理论、人工神经网络(ANN)、支持向量机 (SVM)等^[2,6,7],进一步发展了非线性径流预报模 型.2000年 Tipping 提出了一种新的稀疏概率模型 相关向量机^[8](relevance vector machine,RVM),该 方法用非线性核函数映射到高维空间,在高维空间 进行线性回归,成功实现非线性向线性转化,同时 基于贝叶斯理论定义模型参数,不仅可以定量预 报,而且能够以概率分布的形式描述水文预报不确 定度,可为水库调度决策分析提供更多的可利用信 息.目前,其已应用到图像分析^[9,10]、信道均衡^[11]等 分类与回归问题,获得了较好的应用效果.

综上所述,确定自身前期影响因子和建立预 报模型,是时间序列分析预测方法的关键,本文首 先对径流时间序列进行相空间重构,挖掘水文系 统多变量因子;然后利用重构后的时间序列建立

收稿日期: 2010-01-20; 修回日期: 2011-10-11.

基金项目:水利部公益性行业专项资助项目(201001024);国家自然科学基金资助项目(51109025).

作者简介: 仕玉治(1980-),男,博士,工程师,E-mail:syz101066@163.com;彭 勇*(1979-),男,博士,E-mail:pyongcuidi@163.com.

RVM 非线性径流预报模型,并采用粒子群优化 (PSO)算法^[12]辨识模型参数;最后应用实例验证 本文模型的有效性.

1 混沌时间序列相空间重构

相空间重构是混沌理论的基础,依据 Takens 理论^[13],对某一混沌时间序列{ $x_i:i = 1,2,...,$ n},只要适当选取时间滞时 τ 和嵌入维数 m,且嵌 入维数满足 $m \ge 2D+1$,其中 D 为饱和关联维数, 即可重构与原未知动力系统具有相同几何特征的 m 维相空间,则相空间中的点可以表示为 $X_i =$ ($x_i = x_{i+\tau} = x_{i+2\tau} = ... = x_{i+(m-1)\tau}$)^T,由 N 个相点组 成的延迟状态向量表示为{ $X_i:i = 1,2,...,N$ },其 中 $N = n - (m-1)\tau$,则相应关联积分表达式为

$$C(r,m) = \frac{2}{(N+1)N} \sum_{i=1}^{N} \sum_{j=i+1}^{N} H(r - \|\mathbf{X}_{i} - \mathbf{X}_{j}\|)$$
(1)

式中: $H(x) = \begin{cases} 0; x \leq 0 \\ 1; x > 0 \end{cases}$, r 为标度尺度, $\| \cdot \| \end{pmatrix}$ 欧几里得范数.

对于混沌时间序列,关联积分 C(r,m) 与标 度尺度 r 近似成指数关系: $C(r,m) \propto r^{D}$,

$$D = \lim_{r \to 0} \lim_{m \to \infty} \frac{\partial \ln C(r,m)}{\partial \ln r}$$
(2)

2 相关向量机(RVM)径流预报模型

已知相关向量{ $X_i: i = 1, 2, \dots, N$ },给定任意 一个输入向量 X^* ,则通过非线性映射到高维特 征空间,然后在高维特征空间中进行线性回归得 到预报输出值 \hat{y} ,即 RVM 的模型输出表示为

$$\hat{y} = \sum_{i=1}^{N} w_i K(\boldsymbol{X}^*, \boldsymbol{X}_i) + w_0 \qquad (3)$$

式中:K(•,•)为核函数,w为模型的权值.

令训练样本集为 $\{X_i, y_i\}_{i=1}^{N}$,且p(y | X)服从 $N(y | \hat{y}, \sigma^2)$ 的高斯分布,则相应的训练样本集的 高斯似然函数为

$$p(\mathbf{y} \mid \mathbf{w}, \sigma^2) = (2\pi\sigma^2)^{-N/2} \exp\left\{\frac{-\|\mathbf{y} - \mathbf{\Phi}\mathbf{w}\|^2}{2\sigma^2}\right\} \quad (4)$$

式中: $\mathbf{y} = (y_1 \quad y_2 \quad \cdots \quad y_N)$; $\mathbf{w} = (w_0 \quad w_1 \quad \cdots \quad w_N)$; $\boldsymbol{\Phi} \ \exists N \times (N+1)$ 核函数矩阵, $\boldsymbol{\Phi}_m = K(\mathbf{X}_n, \mathbf{X}_n)$, $\boldsymbol{\Phi}_{n1} = 1$. 利用极大似然函数法估计 \mathbf{w}, σ 时会导致严重的过拟合现象^[8], 因此, 为每个 \mathbf{w} 定义高斯先验概率分布函数

$$b(\boldsymbol{w} \mid \boldsymbol{\alpha}) = \prod_{i=0}^{N} N(\boldsymbol{w}_i \mid 0, \boldsymbol{\alpha}_i^{-1})$$
 (5)

然后基于贝叶斯准则计算权值的后验概率分布,即

$$p(\boldsymbol{w} \mid \boldsymbol{y}, \boldsymbol{\alpha}, \sigma^{2}) = (2\pi)^{-(N+1)/2} \mid \boldsymbol{\Sigma} \mid^{1/2} \times \exp\{-(\boldsymbol{w} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\boldsymbol{w} - \boldsymbol{\mu})/2\}$$
(6)

式中:后验分布的协方差和均值分别为 $\Sigma = (\sigma^2 \Phi^T \Phi + A)^{-1}, \mu = \sigma^2 \Sigma \Phi^T y, 其中 A = diag{\alpha_0, \alpha_1, \dots, \alpha_N}.$

通过最大化边缘似然分布函数

 $p(\mathbf{y} \mid \boldsymbol{\alpha}, \sigma^2) = N(0, \sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{A}^{-1} \boldsymbol{\Phi}^{\mathrm{T}}) \quad (7)$ 即可得超参数估计值,本文采用 EM 算法^[8] 内循环 迭代估计超参数 $\boldsymbol{\alpha}, \sigma, \mathbf{x}, \mathbf{\alpha}_i$ 和 σ 的迭代方程分别如下:

$$\alpha_i^{\text{new}} = \frac{1}{\Sigma_{ii} + \mu_i^2} \tag{8}$$

$$(\sigma^2)^{\text{new}} = \| \mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\mu} \|^2 / (N - \sum_i \gamma_i) \qquad (9)$$

式中: $\gamma_i = 1 - \alpha_i \Sigma_{ii}$.

获得超参数最优值 $\alpha_{mp}, \sigma_{mp}^2, 则 w = \mu_{mp}, 由于$ 权值最优估计的不确定性,给定任意一个输入值 $X^*, 则描述预测不确定度的均值和方差如下:$ $\mu^* = \Phi(X^*)\mu_{mp}, \sigma_*^2 = \sigma_{mp}^2 + \Phi^T(X^*) \Sigma \Phi(X^*)$ (10)

3 模型参数辨识及计算流程

RVM 径流预报模型需要解决两个问题:(1) 模型核函数选择;(2)模型参数辨识过程中的目标 函数确立.对于核函数的选择,线性核函数是径向 基核函数的特例,特定的 sigmoid 核函数功能上 与径向基核函数相同,核函数自身参数的个数太 多不利于参数的选择^[7],因此,本文选取径向基核 函数为核函数;其次,在模型参数识别过程中,通 常选取训练样本的拟合误差最小为模型目标函 数,但是该方式下训练误差收敛过程中会出现严 重的过拟合现象,如图1中拟合曲线1所示,训练 阶段拟合误差非常小,几乎接近于零,导致优化参 数不合理,外推期预报精度非常低.因此本文在训 练过程中考虑具有丰、平、枯年份的检验样本误差 来抑制过拟合,即综合考虑训练样本和检验样本 的误差建立目标函数,如下式:

min $R = \left(\frac{N_2}{N_1 + N_2}\right) R_1 + \left(\frac{N_1}{N_1 + N_2}\right) R_2$ (11)

式中: R_1 、 R_2 分别为训练期、检验期的相对误差绝 对值的平均值, N_1 、 N_2 分别为训练期、检验期的 样本个数. 其收敛过程如图 1 中拟合曲线 2 所 示,曲线 2 的收敛值要比曲线 1 的收敛值大,说明 本文目标函数具有抑制过拟合的能力.

最后,利用 RVM 模型进行径流预报的主要步骤如下:

(1)给定参数 *m*、τ、ε 的合理取值区间,对时 间序列进行相空间重构;



图1 不同目标函数的训练误差收敛过程



(2)利用重构后的训练样本作为 RVM 的输入进行训练,采用 PSO 算法辨识方法参数,并验证数据序列的混沌特性,在 RVM 内循环中,利用式(8)、(9)迭代估计超参数 α_i 、 σ ,给定一个 α_{max} 的值,将 $\alpha_i^{new} > \alpha_{max}$ 对应的权重值设定为 0,并剔除对应的 X_i ,所剩余的 X 即为相关向量,对应的权重 向量 $w = \mu_{mp}^{T}$;

(3)依据优化所得参数为模型的参数取值,给 定任意一个输入向量 X^* ,采用训练好的 RVM 进 行计算便可得到预报值的均值 μ^* 和方差 σ^2_* ,预报 值服从均值 μ^* 和方差 σ^2_* 的后验正态概率分布.

4 应用实例分析

选取南方两水库入库月径流时间序列作为研 究实例,分别为水库1的51a(1953-01~2003-12) 和水库 2 的 48 a(1958-01~2005-12)入库月径流时 间序列.水库1的控制流域面积为11.45×104 km²,多年平均月径流量为1251 m³/s,变差系数为 0.717,最大、最小月径流量分别为5000、248 m^3/s ;水库2的控制流域面积为10.26×10⁴ km², 多年平均月径流量为1215 m³/s,变差系数为 0.856,最大、最小月径流量分别为5 480、236 m³/s.每一个样本序列分成3个子样本,对于水库 1,将前41 a 资料作为模型的训练样本,中间包含 丰、平、枯年份的5a资料作为检验样本,与前41 a 配合确定合理的模型参数,剩余5a 不参加确定 模型参数,纯粹用于检验确定模型的外推预报能 力. 同样,对于水库2,用前38 a 的序列点作为训 练样本,中间5a作为检验样本,剩余5a序列点 作为外推预测样本.对数据进行规格化处理,采用 PSO 算法优化各方法中的参数,取相对误差的绝 对值(E_{mar})、相关系数(R)、确定性系数 R² 和合格 率(定量)作为预报结果的评价指标.先以水库1 为例进行计算分析,以水库2作进一步的验证.

4.1 实例分析

RVM 模型参数主要有相空间重构参数(m、 τ)、所选取的核函数自身参数及模型自动确定的 超参数(α , σ).给定一个较小的参数区间,进行优 化计算,若优化所得参数值为区间端点值,则进一 步扩大区间重新计算,直至参数取值在区间范围 内为止,即该区间为参数区间,由此确定径向基核 函数带宽 ε,混沌时间序列嵌入维数 m、时间滞时 τ 的取值区间分别为[0.1,100]、[1,20]、[1,10], 另外对超参数初始化,取 α (0) = (0.25,0.25,..., 0.25), σ_0^2 = var(y)×0.01, α_{max} = 1×10⁵.

其次,对时间序列进行相空间重构,以重构后 的训练样本作为径流模型的输入条件,采用 PSO 优化模型参数(ϵ, m, τ),水库1结果为(2.1723, 14,4),水库2结果为(4.554,7,6).根据优化所得 时间滞时,分别以嵌入维数 $m = 1 \sim 20$ 绘制 D-log₂ r 折线斜率图, 如图 2 所示, log₂ r 在 1~2, 且 m>8 时饱和关联维数趋于稳定值,即存在显 著标度区,从而可以定性判断两水库月径流序列 存在混沌特性,并由图 2(a)、(b)可以估计出 1< D < 3,重构参数 m 满足 m ≥ 2D+1 的条件,说明 重构参数是合理的.此外,在参数内循环过程中, 随着超参数的迭代估计,边缘似然分布函数值逐 步趋于稳定,如图 3 所示;由图 4 知超参数 σ^2 收 敛很快,迭代3次后基本达到最优值,因此有的文 献也将超参数 σ² 作为一个固定值进行内循环计 算.





Fig. 4 Convergence process of hyper-parameter σ^2

分析模型的模拟、检验、外推预报精度,并将 本文模型(RVM)与应用较为广泛的最小二乘支 持向量机模型(LSSVM),以及未考虑相空间重构 的(*m*=12,τ=1)相关向量机模型(RVM*)和自 动回归滑动平均模型(ARMA(5,6))进行对比分 析.计算结果列于表 1 中,RVM 方法对于训练 期、检验期和外推预测期的预报结果见图 5~8.

总体而言,由表1知,考虑相空间重构进行预 报时比未考虑相空间重构时,RVM获得比单一 方法更优的预报精度,与LSSVM和ARMA(5, 6)的计算结果相比较,RVM的评价指标值均优 于其相应值,说明RVM具有较好的预报性能.按 照《水文情报预报规范》(SL 250—2000)标准将径 流量划分为枯、偏枯、平、偏丰、丰5个级别,对高 流量(包括偏丰和丰流量)精度进行了定量、定性 比较分析,结果如表2所示.RVM在训练期、检 验期及外推预测期的平均绝对相对误差分别比 LSSVM和ARMA的相应值要小,但同时其预报 精度均比预报总体时相应值低,以多年变幅的 20%为许可误差,比较分析知,其定量合格率较本

表1 水库1月流量不同方法预测精度

方法	训练期				检验期				外推预测期			
	$E_{ m mar}/\%$	R	合格率/%	R^2	$E_{ m mar}/\%$	R	合格率/%	R^2	$E_{ m mar}/\%$	R	合格率/%	R^2
RVM	12.34	0.94	82.2	0.88	11.92	0.95	83.3	0.91	13.59	0.93	76.7	0.81
RVM*	15.62	0.92	75.4	0.83	17.58	0.89	68.3	0.79	14.93	0.92	66.7	0.75
LSSVM	13.76	0.93	76.1	0.86	13.37	0.93	76.7	0.85	14.31	0.92	73.3	0.77
ARMA(5,6)	18.27	0.91	62.2	0.72	17.75	0.89	61.7	0.72	17.07	0.93	65.0	0.70



Fig. 6 Comparison and scatter plot between observed flow and predicted flow by RVM during test period

Fig. 7 Comparison and scatter plot between observed flow and predicted flow by RVM during validated period

图 8 外推预测期实测流量与 RVM 预报区 间对比图

文所列其他方法有所提高.为提供更为充分的预 报信息,本文对比分析了三阶段高流量的定性预 报合格率,除了在检验期 RVM 和 LSSVM 的合 格率相同以外,其余两阶段 RVM 均获得比其他 方法更高的定性预报合格率,同样说明 RVM 具 有较强的高流量预报能力.

进一步考虑径流预报的不确定性,以预报值 的均值和方差为预报的后验概率分布函数来描述 预报值的不确定性,并讨论了发生概率为 80%的 区间预报,其区间预报结果及实测流量过程如图 8 所示.由图 8 知,中低流量预报区间基本上可以 包住实测流量,高流量区间上下限值对应的级别 能够预报出实测值对应的级别,概率区间预报是 可靠的.

4.2 实例验证分析

水库 2 的统计参数与水库 1 基本相同,但是 变差系数较大,数据序列平稳性相对较差,在同样 可行条件下,对水库 2 进行了计算,其预报结果的 评价指标列于表 3.由表 3 知, RVM 较 LSSVM 和 ARMA(6,6)模型具有较高的预报精度,验证 说明了本文模型的有效性.

表 2 水库 1 高月流量不同方法预测精度

Tab. 2 Prediction accuracy of high monthly flow of Reservoir One resulting from various methods

方法 -		训练期			检验期		外推预测期			
	$E_{ m mar}/\%$	定量合格率/%	定性合格率/%	$E_{ m mar}/\%$	定量合格率/%	6定性合格率/%	$E_{ m mar}/\%$	定量合格率/	'% 定性合格率/%	
RVM	14.27	75.4	92.9	14.92	73.9	82.6	16.53	70.8	83.3	
LSSVM	18.24	71.6	87.8	18.94	67.8	82.6	18.13	66.6	79.1	
ARMA(5,6)	19.03	62.1	74.3	22.41	56.5	70.0	19.32	62.5	76.6	

表 3 水库 2 月流量不同方法预测精度

Tab. 3	Prediction	accuracy of	monthly flow	of Reservoir	Two 1	resulting	from	various	methods
--------	------------	-------------	--------------	--------------	-------	-----------	------	---------	---------

方法	训练期				检验期				外推预测期			
	$E_{ m mar}/\%$	R	合格率/%	R^2	$E_{ m mar}/\%$	R	合格率/%	R^2	$E_{ m mar}/\%$	R	合格率/%	R^2
RVM	13.00	0.938	81.7	0.86	11.17	0.959	80.0	0.89	14.62	0.907	75.0	0.84
LSSVM	15.68	0.921	76.6	0.85	12.64	0.943	73.3	0.79	15.65	0.901	71.7	0.76
ARMA(6,6)	18.74	0.903	60.7	0.71	18.98	0.900	56.7	0.71	19.15	0.885	51.7	0.67

5 结 论

(1)将混沌技术与相关向量机结合建立径流 预报模型,采用 PSO 算法辨识模型参数,优化所 得重构参数满足混沌理论条件,耦合方法比单一 方法的预报精度有所提高,并对总体和高流量值进行分析,取得比LSSVM和ARMA模型更优的预报精度,说明本文模型的有效性.

(2)相关向量机为概率模型,能够定量地、以 概率分布的形式描述径流预报不确定性,并给出

Fig. 8 Comparison between observed and predicted interval hydrograph during validated period

指定发生概率下的区间预报.

(3)在进行中长期径流预报应用时,相关向量 机模型的不足之处是模型参数和样本序列均以正 态概率分布函数进行推理,但从模型计算的结果 来看可用于中长期径流预报,下一步将以 P-Ⅲ型 概率分布函数进行模型研究.

参考文献:

- [1] SIVAKUMAR B. Chaos theory in hydrology important issues and interpretations [J]. Journal of Hydrology, 2000, 227:1-20
- [2] SIVAKUMAR B, JAYAWARDENA A W, FERNANDO T M K G. River flow forecasting: use of phase-space reconstruction and artificial neural networks approaches [J]. Journal of Hydrology, 2002, 265:225-245
- [3] 王 文,许武成.对水文时间序列混沌特征参数估计 问题的讨论[J].水科学进展,2005,16(4):606-610
- [4] YU X Y, LIONG S Y, BABOVIC V. EC-SVM approach for real time hydrologic forecasting [J]. Journal of Hydroinformatics, 2004, 6(3):209-223
- [5] SIVAKUMAR B. Nonlinear determinism in river flow prediction as a possible indicator [J]. Earth Surface Processes and Landforms, 2007, 32:969-979
- [6] LIONG S Y, SIVAPRAGASAM C. Flood stage forecasting with SVM [J]. Journal of the American

Water Resources Association, 2002, 38(1):173-186

- [7] 林剑艺,程春田. 支持向量机在中长期径流预报中的 应用[J]. 水利学报,2006,**37**(6):681-686
- [8] TIPPING M E. The relevance vector machine [J]. Advances in Neural Information Processing System, 2000, 12:652-658
- [9] AGARWAL A, TRIGGS B. 3D human pose from Silhouettes by relevance vector regression [J]. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004, 2:882-888
- BOWD C, MEDEIROS F A, ZHANG Zuo-hua, et al. Relevance vector machine and support vector machine classifier analysis of scanning laser polarimetry retinal nerve fiber layer measurements
 Investigative Ophthalmology & Visual Science, 2005, 46:1322-1329
- [11] CHEN S, GUNN S R, HARRIS C J. The relevance vector machine technique for channel equalization application [J]. IEEE Transactions on Neural Networks, 2002, 12(6):1529-1532
- [12] KENNEDY J, EBERHART R C. Particle swarm optimization [C] // Proceedings of IEEE Conference on Neural Networks. Piscataway: IEEE Press, 1995: 1942-1948
- [13] KANTZ H, SCHREIBER T. Nonlinear Time Series Analysis [M]. Cambridge: Cambridge University Press, 1997

Research on mid- and long-term runoff forecast model with relevance vector machine

SHI Yu-zhi^{1,2}, PENG Yong^{*1}, ZHOU Hui-cheng¹

(1. School of Hydraulic Engineering, Dalian University of Technology, Dalian 116024, China;
2. Water Research Institute of Shandong Province, Jinan 250013, China)

Abstract: Due to the superior forecasting performance, relevance vector machine (RVM) was applied to mid- and long-term runoff forecasting, and based on the phase space reconstruction, the runoff relevance vector machine forecasting model was established. Firstly, the runoff time series was reconstructed in the phase space, and the reconstructed series was as the proposed model input; Secondly, the particles swarm optimization (PSO) algorithm was applied to identifying the model parameters and chaotic properties of time series. The EM algorithm was used to estimate hyperparameters in the inherent cycle, RVM was compared with widely used least squares support vector machine (LSSVM) and auto-regressive moving average model (ARMA). The test results show that the model has good generalization ability; Finally, in terms of the uncertainty of hydrological processes, the scheduling staffs consider the uncertainties in forecasting, and quantitatively estimate the risks and benefits in decision-making based on the uncertainty of RVM output values and the probability forecast interval.

Key words: phase-space reconstruction; relevance vector machine; long-term runoff forecast; PSO algorithm