

DOI:10.19853/j.zgjsps.1000-4602.2025.07.009

# 基于滤波重构时间序列回归算法的需水量预测

关思源, 张巧珍

(平行数字科技<江苏>有限公司, 江苏 苏州 215200)

**摘要:** 为合理分配供水量、提高供水效率、保障供水安全,自来水厂对用户需水量的准确预测是十分必要的。采用基于样本特征的滤波重构自回归模型,能够在保留趋势性数据的同时,对异常数据进行剔除。通过相关性分析发现,城市需水量与时间具有高度的线性相关性;采用滑动窗口对原数据进行时间序列分析,结合滤波重构,平均绝对百分比误差为2.26%,且近92%的数据落在误差5%以内,明显优于单一时序分析法和机器学习。采用该算法进行需水量预测,建模后应用于某自来水公司生产调度,其降低了出厂水压力和供水压力电耗,降低幅度分别为5.64%和4.55%。

**关键词:** 需水量预测; Savitzky-Golay滤波器; 相关性分析; 滤波重构; 时间序列

**中图分类号:** TU991 **文献标识码:** A **文章编号:** 1000-4602(2025)07-0063-06

## Water Demand Prediction Utilizing Filtering Reconstruction Time Series Regression Algorithm

GUAN Si-yuan, ZHANG Qiao-zhen

(Avatar Digital Technology <Jiangsu> Co. Ltd., Suzhou 215200, China)

**Abstract:** Accurate prediction of user water demand is essential for water treatment plants to rationally allocate water supply, improve water supply efficiency, and ensure water supply safety. The auto-regressive model, reconstructed through feature-based filtering, effectively eliminates anomalous data while preserving trend information. Through correlation analysis, it was determined that there existed a significant linear relationship between urban water demand and time. By employing a sliding window approach to analyze the time series of the original data in conjunction with filtering reconstruction, the mean absolute percentage error was reduced to 2.26%, with approximately 92% of the data points exhibiting an error range within 5%. This performance is notably superior to that achieved through single time series analysis and other machine learning methods. The water demand prediction model, developed using the specified algorithm, was implemented in the production scheduling of a water company. As a result, the power consumptions associated with product water pressurization and water supply pressurization were reduced by 5.64% and 4.55%, respectively.

**Key words:** water demand prediction; Savitzky-Golay filter; correlation analysis; filtering reconstruction; time series

目前,城市供水安全已逐渐成为民生关注的焦点,为了提高供水效率、保障供水安全、优化供水能力分配,自来水厂准确预测需水量至关重要。需水量预测受多种因素的影响,包括天气、温度、节假日、

供水用途等,呈现出不确定性和时空变化性。在数据预测领域,复杂的机器学习回归方法已被广泛应用于模型构建,然而,这类算法通常具有计算复杂、对算力要求高等缺点<sup>[1-2]</sup>;另一方面,算法的选择还

需要考虑到运行环境的兼容性。基于上述原因,有必要选择能够兼容较低配置,同时对误差有较高适应性的算法。

线性回归分析和时间序列分析是常见的计算方法,适用于数据质量良好的情况,并能够充分发挥强大的计算效率,然而,当面临数据异常情况时,这些方法往往表现出较低的适应性。在当前的研究中,所选取的数据集呈现出较高的不规则性,主要体现为数据信息不全、上下波动较大等方面,这些数据异常妨碍了多种预测方法的有效运用。在此背景下,有必要采用灵活度更高且鲁棒性更强的技术,比如滤波重构等数据处理方法,以应对数据集的非理想性,同时又能够保留数据的核心特征,从而达到更贴近实际的拟合效果。此类方法能够有效减少异常数据对分析结果的影响,确保从数据中获得准确可靠的结论。滤波重构在解决异常数据问题时具备显著的优势,在异常值的筛选与剔除方面表现出色,对提升训练及拟合曲线的贴合度方面具有重要价值,这有助于实现对未来数据点的精准控制。

笔者通过引入滤波重构技术,并辅助时间序列分析对自来水厂需水量数据进行处理,揭示其中的潜在规律性和趋势性信息,从而实现更为精准的需水量预测,科学制订供水和调度方案,减少设备的频繁启停,从而降低自来水厂的供水电耗。采用该算法对时间序列数据进行分解,将其分解成不同的分量,并运用回归方法对这些分量进行建模,从而提升预测的准确性和可靠性。

## 1 研究方法

### 1.1 相关性分析

为了系统性地分析需水量变化的相关因素,采用了相关性分析方法,以确定与需水量变动密切相关的因素。相关性分析有助于揭示潜在的线性或非线性关系,以便进一步分析和解释不同因素对需水量变化的影响机制。

#### 1.1.1 线性相关

皮尔逊相关系数是一种用于衡量两个连续变量之间线性关系强度和方向的统计量,使用该系数时需要确保数据的连续性,且不存在明显异常值<sup>[3]</sup>,因此,在使用此方法进行分析前,需要进行初步的空值填补和异常值的处理及优化。本研究使用基

于密度的空间聚类(DBSCAN)算法来筛查并删除离群点,然后使用往年数据的均值来填补空值。将数据集中的核心数据与潜在影响因子进行一一配对,得到序列 $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ,  $x$  和  $y$  的皮尔逊相关系数 $r_{xy}$ 为:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

式中: $n$ 为样本量; $\{x_i, y_i\}$ 为样本值; $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ 。

$r_{xy}$  的取值范围在 $-1 \sim 1$ 之间。当 $r_{xy}$ 趋近于1时,该潜在影响因子与数据走势呈强正相关;当 $r_{xy}$ 趋近于 $-1$ 时,则呈强负相关;当 $r_{xy}$ 趋近于0时,该因子影响性较弱。

#### 1.1.2 非线性相关

皮尔逊相关系数对线性关系较为敏感,但对于非线性相关的灵敏度较为有限。为了确认原数据的线性贴合度,仍需要对原数据进行非线性分析。Kendall's Tau-b 相关分析是基于秩次数据的比较,通过比较数据中各观测值的排名顺序,衡量两个变量之间的相关性。该方法首先比较两个变量中的每对观测值,计算出在排名上一致性和不一致性的数量,然后通过这些数量计算出 Tau-b 相关系数。如果两个变量的排序趋势一致,相关系数接近1,如果排序趋势相反,则相关系数接近 $-1$ 。通过对比每对观测值之间的秩次差异,Kendall's Tau-b 分析能够在非参数条件下评估变量之间的关联程度,尤其适用于非线性的相关性分析<sup>[4]</sup>。Kendall's Tau-b 相关性系数( $\tau_B$ )定义如下:

$$\tau_B = \frac{n_c - n_d}{\sqrt{(n_0 - n_1)(n_0 - n_2)}} \quad (2)$$

式中: $n_0 = \frac{n(n-1)}{2}$ ;  $n_1 = \sum_i \frac{t_i(t_i-1)}{2}$ ;  $n_2 = \sum_i \frac{u_i(u_i-1)}{2}$ ;  $n_c$ 为协和对数; $n_d$ 为非协和对数。

#### 1.2 时间序列分析

时间序列分析是一种用于研究呈时间变化数据的统计方法,旨在揭示数据的随时变性,例如总体趋势和周期性。时序数据是按时间顺序排列的数据集,可以是连续的时间点或者固定时间间隔的数据观测值。本研究采用滑动窗口方法进行时序

分析。

滑动窗口是一种在数据处理、时间序列分析和机器学习中广泛应用的方法。其核心是通过移动一个固定大小的窗口,从数据流或序列数据中逐步提取信息<sup>[5-6]</sup>。通过在每个窗口内提取特征,并滑动各个窗口进行训练,生成多个训练样本,可完成对于序列数据的特征提取,获取以历史数据为基础的训练模型。如果在窗口内应用数据统计指标计算,滑动窗口能实现移动平均、趋势分析、周期性探测等任务,并揭示数据的趋势和模式。其中,每个窗口的大小可以根据任务需求设定,以决定捕捉不同时间尺度下的信息变化。滑动窗口伪代码如下:

$$w[i:j] = \{x[i], x[i+1], \dots, x[j-1]\} \quad (3)$$

式中: $w[i:j]$ 集合了 $i \sim j-1$ 索引的数据; $x[i] \sim x[j-1]$ 代表时间序列中每个节点的历史值,每个节点之间的空隙均可进行多次分割。

### 1.3 滤波分析

为了减少数据中的噪声干扰,使用滤波分析进行数据曲线的平滑处理。滤波器设计的目的是搭载在示波器的微处理器中,对仿真信号进行波形降噪,以揭示出信号的真实趋势。鉴于水量数据与波形图的时间序列共性,以及高度类似的频率响应,使用滤波器滤通强相关点、滤阻离群点具有一定的合理性。

$$\begin{pmatrix} x_{t-n} \\ \vdots \\ x_t \\ \vdots \\ x_{t+n} \end{pmatrix} = \begin{pmatrix} 1 & t-n & (t-n)^2 & \cdots & (t-n)^{k-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t & t^2 & \cdots & t^{k-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t+n & (t+n)^2 & \cdots & (t+n)^{k-1} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{k-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_{t-n} \\ \vdots \\ \varepsilon_t \\ \vdots \\ \varepsilon_{t+n} \end{pmatrix} \quad (5)$$

$F$ 为预测值,根据最小二乘法,解得式(6)。

$$F = B(B^T + B)^{-1} + B^T \quad (6)$$

对某段曲线的滤波重构进行可视化,如图1所示(平均误差为0.003)。

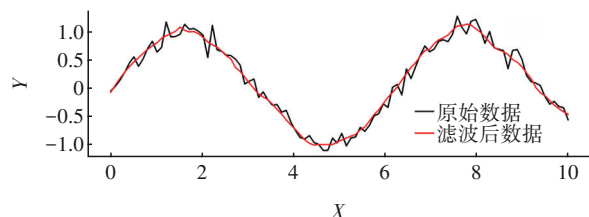


图1 某段曲线的滤波重构可视化

Fig.1 Visualization of filtering reconstruction for a certain curve segment

通过将信号传递给合适的滤波器,选择性地保留或抑制不同频率的成分,异常数据能够得到较为平滑的滤阻。滤波分析的类型多种多样,包括低通滤波、高通滤波、带通滤波等。由于数据低频特征显著丰富,可采用低通滤波器对数据特征进行滤通。滤波分析能够平滑数据、消除噪声、集合趋势以及捕捉周期性成分,从而准确识别有价值的数据范式。

Savitzky-Golay 滤波器是一种基于滑动窗口内的多项式拟合,以此达成对窗口内数据点的预测。该滤波器将时间序列转化为频率序列以后,对其进行窗口化分区,从而达成时间域转换;而后,Savitzky-Golay 滤波器调用高斯拟合函数,依据每个窗口化向量空间的特征方程,对连通变频分量进行线性逼近,依据拟合曲线的空间中心、曲线的卷积,对每个窗口内的数据进行一定的平滑拟合,最后再通过傅里叶分析,将所有分量构成的曲线进行连接<sup>[7]</sup>。Savitzky-Golay 滤波器的平滑功能本质上是一种加权滤波。

对于给定时序点  $2n-1$  个观测值的滤波过程,该滤波器使用  $k-1$  阶多项式进行拟合,见式(4)。

$$x_t = a_0 + a_1 t + a_2 t^2 + \cdots + a_{k-1} t^{k-1} \quad (4)$$

对于  $(t-n, t+n)$  时刻的预测值,构成式(5)所示的矩阵,从左到右的矩阵分别用  $X$ 、 $B$ 、 $A$ 、 $E$  表示。

## 2 案例分析

### 2.1 案例介绍

苏州某自来水公司服务面积为  $1\,176\text{ km}^2$ ,服务规模约 62 万户,总供水能力为  $90 \times 10^4\text{ m}^3/\text{d}$ ,下辖两座水厂,其中第一水厂的处理量为  $60 \times 10^4\text{ m}^3/\text{d}$ ,第二水厂的处理量为  $30 \times 10^4\text{ m}^3/\text{d}$ ,配套供水管网  $7\,500\text{ km}$  (DN75 以上),设置 6 座增压泵站。第一水厂共有 4 组处理单元,每组处理量为  $15 \times 10^4\text{ m}^3/\text{d}$ ,清水库有效容积为  $12 \times 10^4\text{ m}^3$ ,送水泵房安装了 6 台送水泵,单台泵参数:流量  $Q=4\,500\text{ m}^3/\text{h}$ 、扬程  $H=40\text{ m}$ 、电压  $U=6\text{ kV}$ 、功率  $N=600\text{ kW}$ 、变频控制;第二水厂共有 4 组处理单元,每组处理量为  $7.5 \times 10^4\text{ m}^3/\text{d}$ ,清





2.2 误差评价

采用平均绝对百分比误差(MAPE)为主、决定系数( $R^2$ )为辅的方式进行误差分析。MAPE是用于衡量预测值与真实值误差程度的指标,其计算方法见式(7)。

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right| \times 100\% \quad (7)$$

式中: $A_i$ 为实际值; $F_i$ 为预测值。

$R^2$ 是统计量,用于衡量回归模型的拟合优度。它表示因变量方差的一部分能被模型所解释的比例。 $R^2$ 分数的取值范围为0~1,越接近1表示模型对数据的解释能力越好,其按式(8)计算。

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} \quad (8)$$

式中: $y$ 为原始数据; $\hat{y}$ 为预测值; $\bar{y}$ 为原始数据均值。

算法预测结果见表1。相比传统的回归分析以及XGBoost回归、随机森林、支持向量机回归、简单的深度学习等算法,滑动窗口在平均误差、模型贴合度、误差分布等方面皆达到最优。同时,相较单一使用滑窗进行时序分析,经过滤波重构的滑动窗口,其平均误差降低,预测值对于小于5%的误差分布显著稠密。由于高频分量被滤阻,滤波器面对测试集中的突升突降逼近效果较弱,使得最大误差略微偏高。

表1 算法预测结果

Tab.1 Algorithm prediction results

项目	MAPE /%	$R^2$	误差小于5%的百分比/%	误差大于10%的百分比/%	最大误差/%
滤波重构滑窗	2.26	0.891	91.85	0.74	12.79
单一滑窗	2.47	0.862	87.42	0.66	11.96
纯线性回归	3.08	0.804	82.12	1.32	12.28
XGBoost回归	2.93	0.821	82.78	1.32	12.07
随机森林	3.08	0.797	80.79	3.31	12.68
支持向量机回归	2.94	0.817	82.78	1.32	12.25
LSTM	3.06	0.800	81.46	2.65	12.06
GRU	3.06	0.795	80.13	3.31	13.31

2.3 效果评价

根据需水量预测建模,并建立调度方案,按照调度方案进行水厂水量供应分配,以未进行水量预

测的2022年全年的生产数据指标,如平均供水量、出厂水压力、电耗数据,与进行水量预测建模后2023年全年同口径相应指标进行对比,结果如表2所示。

表2 需水量预测前后送水压力和电耗的变化

Tab.2 Change in water supply pressure and electricity consumption before and after water demand prediction

项目	需水量预测前(2022年)				需水量预测后(2023年)			
	供水量/ ( $\text{m}^3 \cdot \text{d}^{-1}$ )	供水 电耗/ $\text{kW} \cdot \text{h}$	出厂 水压 力/ $\text{MPa}$	供水 压力 电耗/ ( $\text{kW} \cdot \text{h} \cdot \text{MPa}^{-1}$ )	供水量/ ( $\text{m}^3 \cdot \text{d}^{-1}$ )	供水 电耗/ $\text{kW} \cdot \text{h}$	出厂 水压 力/ $\text{MPa}$	供水 压力 电耗/ ( $\text{kW} \cdot \text{h} \cdot \text{MPa}^{-1}$ )
1月	504 132	123.5	0.317	389.56	539 210	110.99	0.293	378.82
2月	434 263	120.0	0.305	393.29	464 961	106.88	0.291	367.29
3月	522 239	118.5	0.302	392.30	545 060	110.06	0.298	369.31
4月	524 075	122.3	0.325	376.27	551 611	111.14	0.310	358.52
5月	525 593	120.4	0.307	392.21	565 071	110.86	0.307	361.10
6月	559 164	121.0	0.309	391.55	588 544	110.26	0.307	359.14
7月	575 180	123.6	0.318	388.76	615 840	109.98	0.306	359.41
8月	633 361	124.7	0.337	370.09	661 037	108.86	0.306	355.74
9月	621 664	123.9	0.335	369.82	654 327	110.32	0.308	358.20
10月	593 277	120.5	0.330	365.21	624 329	107.24	0.293	366.02
11月	552 949	122.1	0.329	370.97	568 996	108.78	0.295	368.74
12月	517 209	117.6	0.315	373.46	544 385	107.76	0.293	367.78
均值	546 925	121.7	0.319	381.47	579 908	109.45	0.301	364.13
注: 电耗均以1 000 $\text{m}^3$ 水计。								

从表2可以看出,采用需水量预测建模应用后与预测前进行对比发现,出厂水压力、供水电耗和供水压力电耗均有一定幅度的降低:出厂水压力由0.319 MPa降至0.301 MPa,降幅为5.64%;每1 000  $\text{m}^3$ 供水电耗由121.72  $\text{kW} \cdot \text{h}$ 降至109.45  $\text{kW} \cdot \text{h}$ ,降幅为10.08%;每1 000  $\text{m}^3$ 供水压力电耗由381.47  $\text{kW} \cdot \text{h}/\text{MPa}$ 降至364.13  $\text{kW} \cdot \text{h}/\text{MPa}$ ,降幅为4.55%。

3 结论

① 为了建立一种较为精确的应用型水量预测机制,针对某自来水厂4年的供水量,采用滤波重构方法对结构受损且受严重异常干扰的原始数据进行了系统性平滑处理,以此提升预测模型的稳定性和准确性,可以满足在实际应用中对可靠性和有效性的高要求。

② 相比于当前主流的机器学习回归模型,使用时序分析在误差评价方面达到局部最优;覆盖滤波器后,采用的滤波重构方法在降低误差方面表现

出明显的优越性;同时,经过滤波重构处理的时序回归模型,不仅在低限差内预测值的分布上表现出优势,在整体预测准确性上也显著提升。由此,滤波重构算法为时序分析领域的精确性与可靠性提供了有力支持。

③ 滤波重构作为一种跨领域的研究方法,其核心思想中融合了对干扰分量的滤阻,以及对富含特征值分量的滤通。城市供水量数据常伴随着各种不确定性,因此滤波重构方法对于此类不平衡数据能够呈现出理想的应用潜力。滤波重构能够有效消除数据中的噪声干扰,同时保留重要的特征信息,使其为实现精准和可靠的供水管理规划提供了新颖且可行的解决方案。

④ 采用滑动窗口对原数据进行时间序列分析,结合滤波重构,平均约对百分比误差为2.26%,并且近92%的数据落在误差5%以内,显著优于单一时序分析法和机器学习。

⑤ 采用基于滤波重构时间序列回归算法的水厂需水量预测,建模后应用于水厂的实时生产调度中,降低了出厂水压力和供水压力电耗,降低幅度分别为5.64%和4.55%,具有较好的应用效果。

#### 参考文献:

- [1] 林昱道,陶涛,信昆仑,等. 基于图深度学习的供水管网短期需水量预测研究[J]. 环境工程, 2023, 41(4): 149-153.
- LIN Yudao, TAO Tao, XIN Kunlun, *et al.* Graph deep learning: application on short-term water demand forecasting for water distribution network [J].

Environmental Engineering, 2023, 41(4): 149-153 (in Chinese).

- [2] YOUNES H, ALAMEH M, IBRAHIM A, *et al.* Efficient Algorithms for Embedded Tactile Data Processing[M]. Gistrup: River Publishers, 2020.
- [3] VUKOVIĆ M, LILAND H K, INDAHL U G, *et al.* Extraction of photoluminescence with Pearson correlation coefficient from images of field-installed photovoltaic modules [J]. Journal of Applied Physics, 2023, 133(21):214901.
- [4] GHALIBAF M B. Relationship between Kendall's Tau correlation and mutual information [J]. Revista Colombiana de Estadística, 2020, 43(1):3-20.
- [5] CHAN L S H, CHU A M Y, SO M K P. A moving-window Bayesian network model for assessing systemic risk in financial markets [J]. PLOS ONE, 2023. DOI: 10.1371/journal.pone.0279888.
- [6] TALEBI-KALALEH M, MEI Q P. A mobile sensing framework for bridge modal identification through an inverse problem solution procedure and moving-window time series models[J]. Sensors, 2023, 23(11):5154.
- [7] SCHAFER R W. What is a Savitzky-Golay filter? [J]. IEEE Signal Processing Magazine, 2011, 28(4): 111-117.

作者简介:关思源(2000—),男,江苏苏州人,硕士,工程师,主要从事水厂和污水处理厂的自动化控制及数字仿真等技术工作。

E-mail:sguan22@wisc.edu

收稿日期:2024-08-23

修回日期:2024-10-10

(编辑:任莹莹)

精打细算用好水资源,从严从细管好水资源