

基于特征选择和SSA-LSSVM的短期PM_{2.5}浓度预测

金春梅

贵州大学, 数学与统计学院, 贵州 贵阳

收稿日期: 2022年2月9日; 录用日期: 2022年3月3日; 发布日期: 2022年3月10日

摘要

为了提升PM_{2.5}浓度的预测精度, 考虑到PM_{2.5}浓度受时间序列特征和非线性特征等原因的影响, 导致了时间序列分析模型在预测PM_{2.5}浓度时会存在较大偏差。为此, 提出一种基于特征选择和麻雀搜索算法(Sparrow Search Algorithm, SSA)优化最小二乘支持向量机(Least Squares Support Vector Machine, LSSVM)参数的定量预测模型。首先, 将14个特征变量进行二进制编码, 利用遗传算法结合最小二乘支持向量机对特征变量进行优选, 获取最优特征子集; 利用SSA算法对LSSVM的参数进行优化, 构建SSA-LSSVM的PM_{2.5}浓度预测模型。实验结果表明, 基于遗传算法进行特征选择和麻雀算法优化最小二乘支持向量机参数的模型, 具有明显的预测效果。其中, 该模型的RMSE和MAE分别为10.53和8.01, 预测精度均高于其它模型。

关键词

最小二乘支持向量机, 遗传算法, 麻雀搜索算法, 特征选择, PM_{2.5}浓度

Short-Term PM_{2.5} Concentration Prediction Based on Feature Selection and SSA-LSSVM

Chunmei Jin

School of Mathematics and Statistics, Guizhou University, Guiyang Guizhou

Received: Feb. 9th, 2022; accepted: Mar. 3rd, 2022; published: Mar. 10th, 2022

Abstract

In order to improve the prediction accuracy of PM_{2.5} concentration, considering the time series characteristics of the PM_{2.5} concentration influencing factor dataset and the nonlinear character-

ristics of the data, the time series analysis model still has a large error in predicting $PM_{2.5}$ concentration. To this end, a quantitative prediction model based on feature selection and Sparrow Search Algorithm (SSA) to optimize the parameters of Least Squares Support Vector Machine (LSSVM) is proposed. First, the 14 feature variables are binary coded, and the feature variables are optimized by using the genetic algorithm combined with the Least Squares Support Vector Machine to obtain the optimal feature subset; the SSA algorithm was used to optimize the parameters of the LSSVM, and the $PM_{2.5}$ concentration prediction model of the SSA-LSSVM was constructed. The experimental results show that the model based on the genetic algorithm for feature selection and the sparrow algorithm to optimize the parameters of the Least Squares Support Vector Machine has obvious prediction effect. Among them, the RMSE and MAE of this model are 10.53 and 8.01, respectively, and the prediction accuracy is higher than other models.

Keywords

Least Squares Support Vector Machine, Genetic Algorithm, Sparrow Search Algorithm, Feature Selection, $PM_{2.5}$ Concentration

Copyright © 2022 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

随着经济的快速发展,大气污染也日趋严重。 $PM_{2.5}$ 作为主要的污染物之一,因其粒径小,含有大量的有毒有害物质,不仅在大气中的停留时间长而且输送距离远,致使空气能见度降低,严重影响空气质量和大气环境,甚至危害人体健康[1]。因此, $PM_{2.5}$ 浓度预测工作对于严重污染事件的预警具有重要的价值和意义。

近几年来,随着机器学习和群智能优化算法的兴起和发展,机器学习被广泛应用于构建 $PM_{2.5}$ 浓度预测模型等智能识别领域,并且具有较高的准确率和实用性。现有研究采用历史气象数据或历史污染数据,利用支持向量回归模型[2]、随机森林[3][4]以及LSTM网络[5][6]等单机器学习模型,预测实时 $PM_{2.5}$ 浓度。王勇等[7]利用FFT与LSTM神经网络方法构建 $PM_{2.5}$ 浓度预测模型,开展未来24h的 $PM_{2.5}$ 浓度预测研究。张旭等[8]利用粒子群优化思想和遗传算法的交叉和变异操作相结合,对BP初始权值和阈值进行设定,并对 $PM_{2.5}$ 浓度预测,可以有效避免陷入局部极小,提高收敛速度。李建更等[9]提出了PLS-MSP(Partial Least Square-MSP)模型用于 $PM_{2.5}$ 浓度预测。马天成等[10]采用一种改进型PSO优化的模糊神经网络,预测 $PM_{2.5}$ 颗粒物浓度的变化规律。但是神经网络大多存在训练时间长、收敛速度慢的缺点。支持向量机在解决小样本、非线性以及高维模型中表现出许多特有的优势,并且具有良好的泛化能力[11]。孟凡念等[12]融合集合经验模态分解和样本熵特征提取理论,提出基于粒子群优化最小二乘支持向量机的方法用于滚动轴承故障识别。最小二乘支持向量机(LSSVM)作为新一代机器学习算法,能够避免神经网络过拟合和支持向量机训练耗时长的的问题[13]。然而LSSVM的参数搜索是采用网格搜索法选择模型参数并进行交叉验证,搜索精度和效率低都是此算法的缺陷[14]。LSSVM建模过程中,LSSVM模型中的参数一般是根据经验来设定的,盲目性大和效率低都是需要解决的问题。为了克服LSSVM算法存在的缺陷,需要采用其他算法优化模型参数,得到最优的参数组合。常用的优化算法一般操作步骤繁琐、耗时较长,可能存在一定的盲目性,并且其精度和收敛速度会因计算的问题维度过高而受到影响。

PM_{2.5}作为主要的大气污染物，其来源和形成过程复杂，影响因素众多，具有高度复杂性和非线性的特征。目前大多数PM_{2.5}浓度预测研究，只是基于PM_{2.5}时间序列，很大程度上制约着预测精度。本文将基于空气质量指标数据(PM_{2.5}、PM₁₀、SO₂、NO₂、CO和O₃)和气象数据(温度、相对湿度、降水量、风速和气压等)对下一日PM_{2.5}浓度进行预测，但是数据当中存在着冗余数据，需要对上述多个特征进行特征筛选。因此，本文采用遗传算法对影响PM_{2.5}浓度的特征变量进行优选，再使用麻雀搜索算法(SSA)来优化LSSVM模型参数，避免了一些传统的优化算法在选择参数时存在的问题，获取最佳的惩罚系数和核函数参数组合，提高了PM_{2.5}浓度的预测精度。

2. 相关工作

2.1. 最小二乘支持向量机

最小二乘支持向量机(LSSVM)是在支持向量机(SVM)的基础上，由 Suykens 等提出的机器学习方法。此模型为一个等式约束优化问题，表示为

$$\min J(\omega, e) = \frac{1}{2} \omega^T \omega + \frac{\gamma}{2} \sum_{k=1}^N e_k^2 \tag{1}$$

式(1)中： $\omega = [\omega_1, \dots, \omega_n]^T$ 为权值系数向量； γ 为惩罚系数； e_k 为误差向量； $e_k, 1, 2, \dots, N$ 。根据 Mercer 条件来定义核函数

$$K = (x_k, x_l) = \varphi^T(x_k) \varphi(x_l) \tag{2}$$

式(2)中： $\varphi(\cdot) = [\varphi_1(\cdot), \dots, \varphi_n(\cdot)]^T$ 为非线性映射函数； $l = 1, 2, \dots, N$ 。最终得到LSSVM的非线性模型

$$y(x) = \sum_{k=1}^N \alpha_k K(x, x_k) + b \tag{3}$$

式(3)中： $\alpha_k (k = 1, 2, \dots, n)$ 为拉格朗日乘子， $\alpha_k \in R$ ； b 为偏差向量。

2.2. 基于遗传算法的特征优选

2.2.1. 染色体编码

由于二进制编码的编码、解码操作简单，交叉、变异等遗传操作便于实现，因此采用二进制编码的方式生成GA染色体，见图1所示。每条染色体由三个基因组成，前两个基因分别由长度为10的二进制编码组成，分别表示LSSVM惩罚因子C和核参数g；第三个基因表示特征变量的选取情况，长度为14，代表14个特征量，1表示对应的特征变量已被选取，0表示对应的特征变量未被选取。

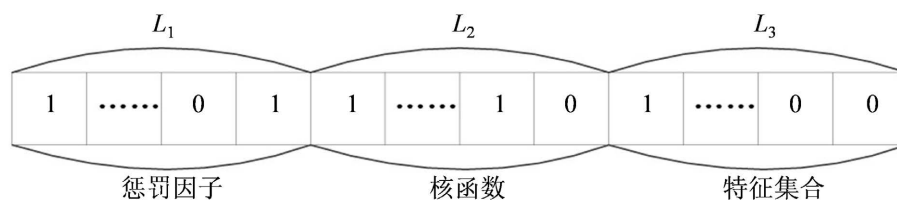


Figure 1. Binary encoding of chromosomes
图 1. 染色体的二进制编码

2.2.2. 适应度计算

对染色体进行编码后，采用训练样本的均方根误差(RMSE)的倒数作为个体适应度 f ：

$$f(L_1, L_2, L_3) = 1 / \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \tag{4}$$

式中 L_1 、 L_2 和 L_3 分别表示最小二乘支持向量机的 C 、 g 和特征变量的组合编码； n 为训练样本数据的数量， \hat{y}_i 为预测结果， y_i 为真实值。

若要得到该染色体的适应度，需要对 L_1 段、 L_2 段和 L_3 段进行分段解码。 L_1 段和 L_2 段解码后得到的十进制数为 LSSVM 的 C 、 g ；根据 L_3 段中每一位的编码值挑选出被选择的特征变量，从而组成新的训练样本。

2.2.3. 遗传操作

遗传操作包括选择操作、交叉操作和变异操作：

- 1) 选择操作能提高了群体的平均适应度值。文中采用“稳态选择”，从而保留父代中适应度较高的个体；
- 2) 交叉操作用于产生新的个体。文中采用“单点交叉”；
- 3) 变异操作用于辅助新个体的产生，其决定了遗传算法的局部搜索能力。文中采用“基本位变异”。

2.2.4. 算法流程

基于 GA 与 LSSVM 的特征变量筛选流程见图 2 所示。

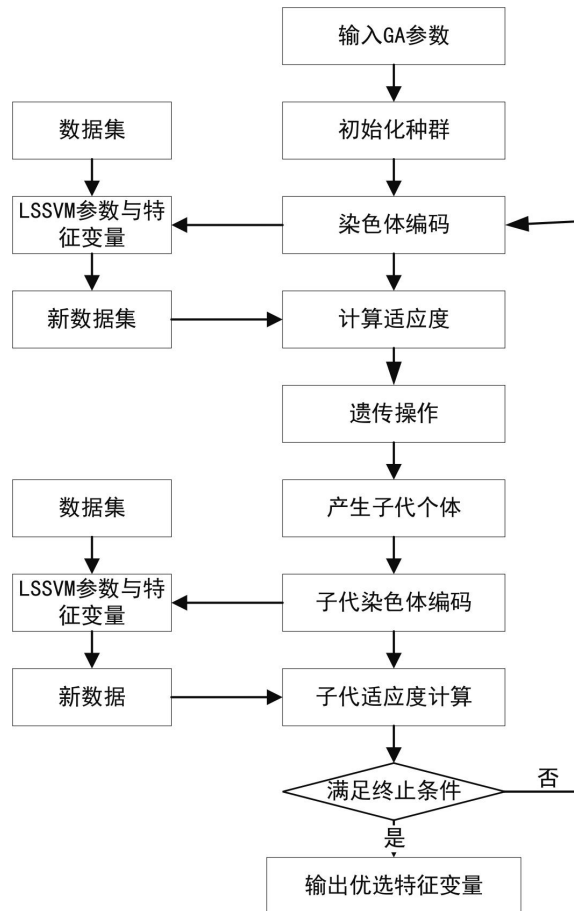


Figure 2. Flow chart of feature selection

图 2. 特征选择的流程图

2.3. 麻雀搜索优化算法

在 SSA 算法中, 每只麻雀代表一个极值优化问题的潜在最优解, 并通过不断更新麻雀种群中发现者、加入者和侦察者的位置, 对目标函数进行优化。SSA 算法中发现者、加入者和侦察者的位置更新公式如下:

$$\mathbf{X}_i^{t+1} = \begin{cases} \mathbf{X}_i^t \times \exp\left(\frac{-i}{\alpha \times t_{\max}}\right) & R2 < ST \\ \mathbf{X}_i^t + q \times \mathbf{L} & R2 \geq ST \end{cases} \quad (5)$$

式中: t 为当前迭代次数, \mathbf{X}_i^{t+1} 表示在第 $t+1$ 次迭代时第 i 只麻雀的适应值, t_{\max} 是最大迭代数, $\alpha \in (0,1)$ 上的一个随机数, $R2$ 表示警戒值, ST 表示安全阈值, q 是一个服从正态分布的随机数, L 是一个一行多维的全一矩阵。

$$\mathbf{X}_i^{t+1} = \begin{cases} q \times \exp\left(\frac{\mathbf{X}_w^t - \mathbf{X}_i^t}{i^2}\right) & i > \frac{n}{2} \\ \mathbf{X}_p^t + |\mathbf{X}_i^t - \mathbf{X}_p^t| \times \mathbf{A}^+ \times \mathbf{L} & i \leq \frac{n}{2} \end{cases} \quad (6)$$

式中: \mathbf{X}_p^t 表示发现者的最佳位置, \mathbf{X}_w^t 表示当前麻雀种群中的最差位置, $\mathbf{A}^+ = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1}$, A 是一个各元素为 1 或 -1 的一行多维矩阵。

$$\mathbf{X}_i^{t+1} = \begin{cases} \mathbf{X}_b^t + \beta \times |\mathbf{X}_i^t - \mathbf{X}_b^t| & f_i \neq f_g \\ \mathbf{X}_i^t + K \times \left(\frac{|\mathbf{X}_i^t - \mathbf{X}_w^t|}{(f_i - f_w) + \varepsilon}\right) & f_i = f_g \end{cases} \quad (7)$$

式中: \mathbf{X}_b^t 表示当前全局最佳位置, β 是步长控制参数, $K \in [-1,1]$ 是一个随机数, f_i 是当前的种群适应度, f_g 和 f_w 分别为当前的最优适应度和最差适应度, ε 表示一个常数, 主要用于避免分母为零, 在本文中 ε 取值为 $10E-8$ 。

2.4. 评价指标

本文中采用常见的评价指标均方误差(MSE)、均方根误差(RMSE)和平均绝对误差(MAE)进行模型精度比较, 预测方法效果更好, 模型预测精度更高。指标的计算公式如下所示:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (8)$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |(y_i - \hat{y}_i)| \quad (9)$$

3. 基于特征选择和 SSA-LSSVM 的 PM_{2.5} 浓度预测

为了解决影响 PM_{2.5} 浓度的特征变量之间存在的相关性或者冗余性问题, 以及最小二乘支持向量机的核参数、惩罚因子难选取等问题。采用遗传算法对特征数据集进行特征选择, 获取最优特征子集。然后, 融合 SSA 算法的快速搜索能力寻找最优的 LSSVM 参数, 构建新的 LSSVM 模型。本文模型的具体操作步骤如下:

- 1) 特征选择。采用遗传算法进行特征选择, 筛选出最优特征子集。

2) 数据预处理。为避免实验数据的取值范围差异和量纲的影响。采用采用最大 - 最小规范化方法对所选特征子集进行归一化处理，公式如下：

$$x^* = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \tag{10}$$

式中： x 表示原始数据的真实值， x^* 表示变换后的数据值， x_{\max} 表示原始数据中的最大值， x_{\min} 表示原始数据中的最小值。

3) 参数初始化。输入最优特征子集、SSA 算法参数。算法参数包括麻雀种群规模 N' ，最大迭代次数 t_{\max} ，搜索空间维数 D ，参数搜索范围以及安全阈值 ST 等。

4) 初始化麻雀位置。基于惩罚因子 C 和核参数 g 等相关参数所设定的搜索范围，采用随机数法初始化麻雀位置 $X_i = (C_i, g_i) (i = 1, 2, \dots, N')$ 。

5) 计算初始种群中每只麻雀对应的适应度值，并确定当前最优适应度 f_g 和最差适应度值 f_w 所对应的位置 X_b 和 X_w 。其中，适应度函数定义为训练样本集的均方误差。

6) 更新位置。根据式(5)~式(7)，依次更新发现者、加入者和侦察者的位置信息。

7) 判断算法是否达到最大迭代次数 t_{\max} ，若是，则循环结束，输出寻优结果；否则返回步骤 6)。

构建融合特征选择和麻雀搜索算法优化最小二乘支持向量机，其具体实现流程见图 3 所示。

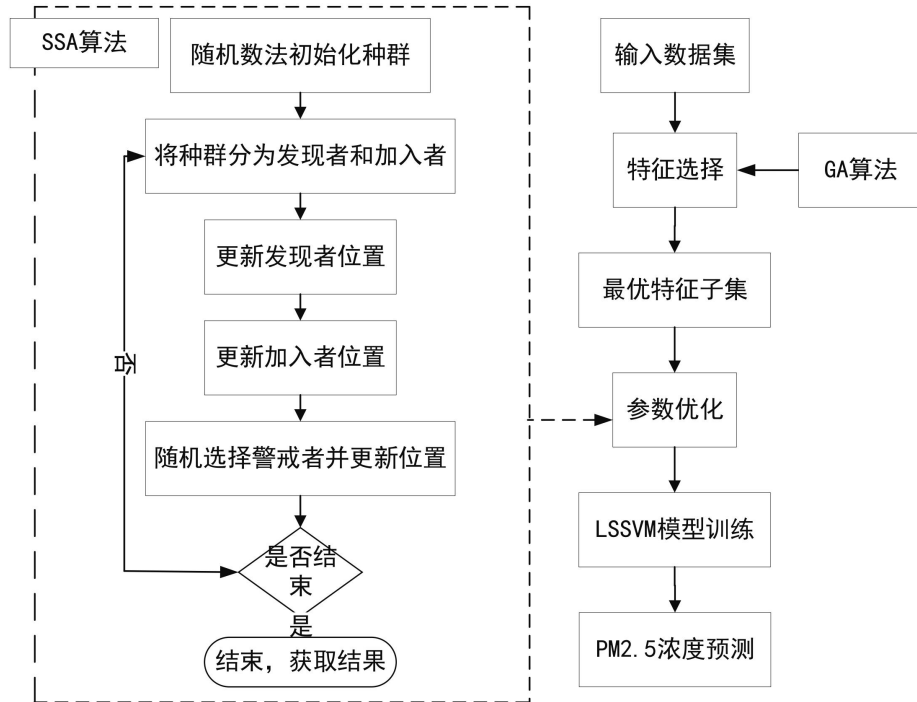


Figure 3. Flow chart of algorithm
图 3. 算法流程图

4. 实验验证

4.1. 数据来源及其预处理

本文中以贵州省贵阳市为例，实验数据包含贵阳市空气质量数据和气象数据，其中空气质量数据来源于空气质量在线监测分析平台，气象数据来源于全国温室数据系统。实验数据选取 2019 年 1 月 1 日

~2020年12月31日的空气质量数据和气象数据的日均值,共有730例样本,并将2019年1月1日~2020年10月31日的日均值数据划分为训练集,2020年11月1日~2020年12月31日的日均值数据划分为测试集。实验数据中收集贵阳市2019~2020年的六个空气质量指标($PM_{2.5}$ 、 PM_{10} 、 SO_2 、 NO_2 、 CO 、 O_3)和八个气象指标(蒸发量、气压、日照时数、降水量、气温、风速、地表气温、相对湿度)作为 $PM_{2.5}$ 浓度的影响因素,并以下一日的 $PM_{2.5}$ 浓度数据作为因变量。为了方便进行特征选择,将变量进行编码见表1。

Table 1. Variable coding

表 1. 变量编码

编码	特征变量	编码	特征变量
1	$PM_{2.5}$	8	气压
2	PM_{10}	9	日照时数
3	SO_2	10	降水量
4	NO_2	11	气温
5	CO	12	风速
6	O_3	13	地表气温
7	蒸发量	14	相对湿度

由于 $PM_{2.5}$ 浓度影响因素的量级有所不同,使用前还需归一化处理数据样本,以消除实用数据的取值范围差异和量纲的影响,所以本文中对实验数据进行归一化处理,以获得更好的预测效果。采用最大-最小规范化,即对原始数据进行线性变换,使得原始数据转换到[0, 1]范围,其中线性变换公式如公式(11)。

在使用GA同时优化LSSVM参数和筛选特征变量时,交叉概率为0.9,变异概率为0.01,种群大小为50、最大迭代次数200。LSSVM参数设置如下:惩罚因子 C 和核参数 g 的范围分别设为[0, 200]和[0, 100],特征变量的个数为14。在本文中,为验证SSA-LSSVM模型的有效性,在 $PM_{2.5}$ 浓度数据集进行测试。实验环境:存4GB,64位的Windows10操作系统,实验工具为Python。在实验中,麻雀搜索算法的参数设置为:种群数量,即每次迭代过程中优化问题的潜在最优解的个数为20;最大迭代次数为50;维度,即优化问题中的寻找的相关参数个数为2,文中指径向基核函数参数 C 和惩罚因子 g ;搜索范围为分别设为[0, 200]和[0, 100];安全阈值ST为0.6。

4.2. 特征变量筛选结果分析

采用GA优选特征子集时,具体的特征量选取情况见表2所示。本文将所选最优特征子集作为最小二乘支持向量机的输入变量,并利用遗传算法和麻雀搜索算法LSSVM模型的参数进行优化处理,得到新的LSSVM预测模型,应用于 $PM_{2.5}$ 浓度预测中。

Table 2. Variable coding

表 2. 变量编码

编码	特征变量	编码	特征变量
1	$PM_{2.5}$	10	降水量
5	CO	11	气温
7	蒸发量	12	风速
8	气压		

4.3. 预测方法对比

为了验证本文算法的有效性,基于相同的最优特征组合,采用基于遗传算法(GA)、粒子群优化算法(PSO)和麻雀搜索算法(SSA)构建的 GA-LSSVM 模型、PSO-LSSVM 模型和 SSA-LSSVM 模型与标准 LSSVM 分别进行 $PM_{2.5}$ 浓度的预测分析,其中,在测试集中,三个模型的预测结果具体见图 4 所示,不同模型的评价指标结果见表 3 所示。

利用以上种方法对 $PM_{2.5}$ 浓度建立定量预测模型。由表 3 的分析结果可以看出,在 LSSVM 模型中, RMSE 和 MAE 分别为 16.06 和 13.04,说明预测曲线的精确度可以进一步提高。应用遗传算法、粒子群优化算法和麻雀搜索算法对 LSSVM 进行优化后,获得了最佳惩罚系数 C 和核函数参数 g 组合,并将其用于预测模型中。用 GA-LSSVM、PSO-LSSVM 和 SSA-LSSVM 预测 $PM_{2.5}$ 浓度变化曲线,拟合效果明显提升,相比于标准 LSSVM 模型而言,在 GA-LSSVM 和 PSO-LSSVM 模型中 RMSE 和 MSE 分别降低到了 15.63 和 12.46、12.75 和 9.85,而在 SSA-LSSVM 模型中 RMSE 和 MAE 分别降低到 10.53 和 8.01,精确度得到了一定程度的提高。另外,将 GA-LSSVM 模型和 PSO-LSSVM 模型与 SSA-LSSVM 模型进行对比发现,后者的预测效果更为显著,相比于前三个模型的 RMSE 和 MAE 都有所下降。通过实验验证表明,采用遗传算法对影响 $PM_{2.5}$ 浓度的特征变量进行特征选择,并结合 SSA-LSSVM 方法来提高 $PM_{2.5}$ 浓度的预测精度是可行的。

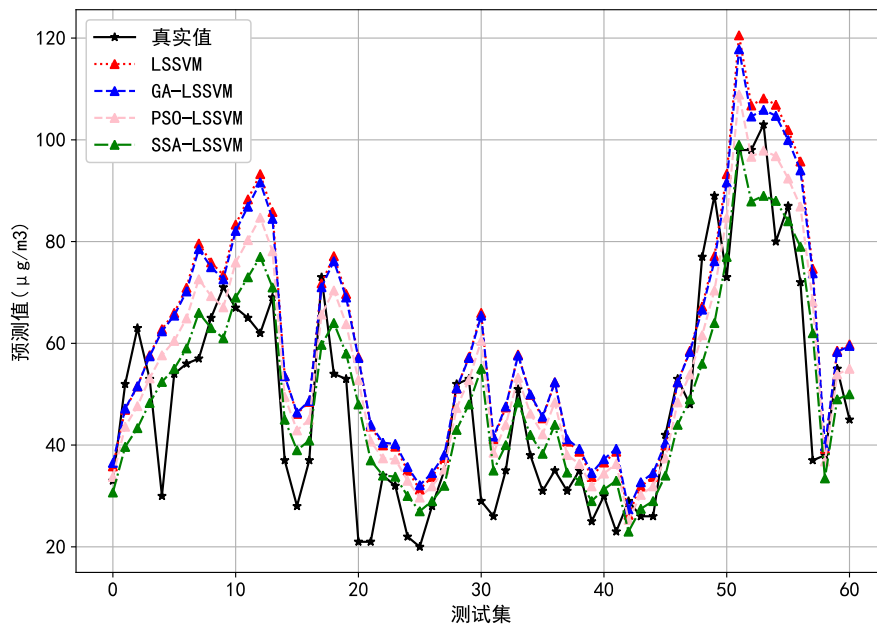


Figure 4. Comparison chart of predicted value and true value

图 4. 预测值与真实值对比图

Table 3. Comparison of evaluation indicators

表 3. 评价指标对比

方法	RMSE	MAE
LSSVM	16.06	13.04
GA-LSSVM	15.63	12.75
PSO-LSSVM	12.46	9.85
SSA-LSSVM	10.53	8.01

5. 结束语

利用遗传算法对影响 $PM_{2.5}$ 浓度的特征变量进行特征筛选, 获取最优特征子集, 用来定量描述 $PM_{2.5}$ 浓度, 将最优特征子集输入 LSSVM、GA-LSSVM、PSO-LSSVM 和 SSA-LSSVM 中进行 $PM_{2.5}$ 浓度的预测。并对结果进行分析, 得出如下结论:

1) 特征变量的选取与 LSSVM 中的惩罚因子和核函数参数对 $PM_{2.5}$ 浓度预测影响都比较大, 对其进行优化能够明显提高预测精度;

2) LSSVM 算法的预测精度依赖于惩罚因子和核函数参数的合理选择, 运用 SSA 算法获取 LSSVM 的最优惩罚因子和核函数参数, 避免了参数选取的随机性, 明显提高预测精度;

3) 提出的方法对 $PM_{2.5}$ 浓度预测的效果比较好, 其中, 本文模型的 RMSE 和 MSE 分别达到 10.53 和 8.01。

参考文献

- [1] 康俊锋, 谭建林, 方雷, 等. XGBoost-LSTM 变权组合模型支持下的短期 $PM_{2.5}$ 浓度预测——以上海为例[J/OL]. 中国环境科学: 2021: 1-16.
- [2] 谢永华, 张鸣敏, 杨乐, 等. 基于支持向量机回归的城市 $PM_{2.5}$ 浓度预测[J]. 计算机工程与设计, 2015, 36(11): 3106-3111.
- [3] 任才溶, 谢刚. 基于随机森林和气象参数的 $PM_{2.5}$ 浓度等级预测[J]. 计算机工程与应用, 2019, 55(2): 213-220.
- [4] 夏晓圣, 陈菁菁, 王佳佳, 等. 基于随机森林模型的中国 $PM_{2.5}$ 浓度影响因素分析[J]. 环境科学, 2020, 41(5): 2057-2065.
- [5] 白盛楠, 申晓留. 基于 LSTM 循环神经网络的 $PM_{2.5}$ 预测[J]. 计算机应用与软件, 2019, 36(1): 67-70.
- [6] 赵文芳, 林润生, 唐伟, 等. 基于深度学习的 $PM_{2.5}$ 短期预测模型[J]. 南京师大学报(自然科学版), 2019, 42(3): 32-41.
- [7] 王勇, 王泓易, 刘严萍, 等. 融合 GNSS 水汽、风速与大气污染物的河北省冬季 $PM_{2.5}$ 浓度预测研究[J]. 大地测量与地球动力学, 2020, 40(11): 1145-1152. <https://doi.org/10.14075/j.jgg.2020.11.009>
- [8] 张旭, 杜景林. 改进 PSO-GA-BP 的 $PM_{2.5}$ 浓度预测[J]. 计算机工程与设计, 2019, 40(6): 1718-1723. <https://doi.org/10.16208/j.issn1000-7024.2019.06.038>
- [9] 李建更, 吴水生. 基于 PLS-M5P 模型的 $PM_{2.5}$ 浓度预测[J]. 计算机与应用化学, 2018, 35(12): 959-970. <https://doi.org/10.16866/j.com.app.chem201812001>
- [10] 马天成, 刘大铭, 李雪洁, 等. 基于改进型 PSO 的模糊神经网络 $PM_{2.5}$ 浓度预测[J]. 计算机工程与设计, 2014, 35(9): 3258-3262. <https://doi.org/10.16208/j.issn1000-7024.2014.09.066>
- [11] 王文涛. 深度学习结合支持向量机在人脸表情识别中的应用研究[D]: [硕士学位论文]. 西安: 长安大学, 2016.
- [12] 孟凡念, 杜文辽, 巩晓赞, 等. 基于粒子群优化最小二乘支持向量机的滚动轴承故障识别[J]. 轴承, 2020(12): 43-50.
- [13] Yang, L., Yang, S., Li, S., et al. (2015) Coupled Compressed Sensing Inspired Sparse Spatial-Spectral LSSVM for Hyperspectral Image Classification. *Knowledge-Based Systems*, **79**, 80-89. <https://doi.org/10.1016/j.knosys.2015.01.006>
- [14] Zenendeboudi, A. (2016) Implementation of GA-LSSVM Modelling Approach for Estimating the Performance of Solid Desiccant Wheels. *Energy Conversion and Management*, **127**, 245-255. <https://doi.org/10.1016/j.enconman.2016.08.070>