

# 基于分解 - 集成方法的铁路客运量预测

刘建熙, 赵依虹, 梁美丽

广东外语外贸大学数学与统计学院, 广东 广州

收稿日期: 2022年11月12日; 录用日期: 2022年12月6日; 发布日期: 2022年12月14日

## 摘要

铁路客运量预测是铁路运输组织管理工作的重要基础和主要依据之一。本文建立多种分解 - 集成方法对全国铁路月度客运量进行预测分析。分别利用集合经验模态分解(EEMD)、奇异谱分解(SSA)和小波分解(WT)将原始序列进行分解, 再分别使用季节差分移动自回归模型(SARIMA)模型和反向传播神经网络(BP)模型及其组合模型对分解后的子序列进行拟合、预测和集成。对比研究发现采用分解 - 集成方法有助于提高相关模型的预测准确性, 且EEMD-SARIMA-BP组合模型在所有模型中预测效果最佳。

## 关键词

预测, 铁路客运量, 集合经验模态分解, 奇异谱分解, 小波变换, BP神经网络

# Railway Passenger Transportation Volume Prediction Models Based on Decomposition-Aggregation Methods

Jianxi Liu, Yihong Zhao, Meili Liang

School of Mathematics and Statistics, Guangdong University of Foreign Studies, Guangzhou Guangdong

Received: Nov. 12<sup>th</sup>, 2022; accepted: Dec. 6<sup>th</sup>, 2022; published: Dec. 14<sup>th</sup>, 2022

## Abstract

The forecast of railroad passenger volume is one of the important foundations and main bases of railroad transportation organization and management. In this paper, we establish multiple decomposition-aggregation methods to forecast and analyze the monthly passenger volume of national railroads. We use the ensemble empirical modal decomposition (EEMD), the singular spectrum decomposition (SSA) and the wavelet analysis (WT) respectively to decompose the original

data series into several sub-series, then we process the forecast by fitting, forecasting and aggregating by the seasonal difference moving autoregressive model (SARIMA) model and back propagation neural network (BP) and their combined model of the sub-series, respectively. We find that the use of decomposition-aggregation methods helps to improve the prediction accuracy, and the combined EEMD-SARIMA-BP model has the best prediction effect among all models.

## Keywords

Forecast, Railroad Passenger Volume, Ensemble Empirical Model Decomposition, Singular Spectrum Decomposition, Wavelet Transform, BP Neural Network

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

铁路是国家重要交通基础设施及重大民生工程,关系着千千万万户的出行和国家经济发展[1]。目前正处于我国铁路的高速发展期间,高速铁路服务网络覆盖范围不断扩大,越来越多人享受到安全快捷的高速铁路服务,高铁出行逐渐成为人们的长短途出行首选[2]。铁路客运量作为反映铁路旅客运输的基本产量指标,可以直接反映铁路客运市场,体现铁路网络的客运服务规模。科学的铁路旅客运输量预测方法,可以建立高速铁路的需求分析,并结合建设标准和财政补贴,完善当前铁路发展水平,最终实现运营模式的盈利和更长足的发展[3][4]。

目前,我国许多学者对于铁路客运量预测的研究主要有传统计量和统计方法、人工智能方法和组合预测方法[4]-[16]。传统的计量和统计方法是基于问题的历史数据并利用模式识别,参数估计,模型检测来建立问题的数学模型。王雷[6]等人根据铁路月客运量的趋势性和周期性,分别采用季节性指数平滑法和季节差分自回归移动平均法建立模型。缪巧芬[7]等人先用 SARIMA 模型对铁路客运量进行建模预测,再同时采用 X-13A-S 季节调整方法重新建模预测,得到了 X-13A-S 季节调整方法较优的结论。人工智能相关的预测方法也已取得了较好的预测效果,如将神经网络理论、灰色系统理论、支持向量机理论、遗传算法理论引入到客运量预测中[8]。侯福均和吴祈宗[4]采用反向传播(BP)神经网络算法研究了铁路客运市场的时间序列。王卓[9]等人对客运量预测的 BP 神经网络模型进行优化,得到了改进后的 BP 神经网络模型的预测效果较标准 BP 模型较优的结论。汪健雄等人[10]为克服 BP 神经网络易计算效率和泛化能力低的问题,结合 Gram-Schmidt 正交化定理,提出了基于双层次 BP 神经网络模型的铁路客运量预测模型。针对客运量时间序列的特性和影响因素,许多学者也提出了新的组合预测方式来提高预测的精度。侯丽敏和马国峰[13]在灰色预测理论的基础上,建立了 GM(1,1)与线性回归的组合模型。贺晓霞[14]等人为了充分考虑客运量的周期变动性,结合灰色预测理论 GM(1,1)与周期扩展模型,得到了与实际值较吻合的预测值。刘琳玥[15]提出了利用主成分分析法去除和减少影响原始铁路客运量因素之间的相关性,再进行 PCA-BP 神经网络模型分析和预测,得到较好的预测效果。

近十几年来,我国学者就铁路客运量预测展开了多维度、多角度的探讨与研究,取得了许多宝贵的成果。但是由于铁路客运量本身季节性较强及外部影响因素较多,许多方法仍存在不足和缺陷。在此基础上,进一步探讨铁路旅客运输量预测模式,探索更多方法的结合是否能为预测带来更高的准确性,本文首先采用季节时间序列模型(SARIMA)和 BP 神经网络模型对铁路月客运量进行预测,再采用集合经

验模态分解(EEMD)、奇异谱分解(SSA)和小波分解(WT)三种分解方法对月客运量数据进行分解, 结合 SARIMA 模型和 BP 神经网络模型进行预测。最后, 通过预测结果比较上述方法的有效性与精确性。

## 2. 季节性差分自回归移动平均模型(SARIMA)

### 2.1. 差分自回归移动平均模型基本原理(ARIMA)

具有下列结构的模型被称作差分自回归移动平均(Autoregressive Integrated Moving Average)模型, 简称为 ARIMA( $p, d, q$ )模型[17]:

$$\begin{cases} \Phi(B)\nabla^d x_t = \Theta(B)\varepsilon_t \\ E(\varepsilon_t) = 0, \text{Var}(\varepsilon_t) = \sigma_\varepsilon^2, E(\varepsilon_t \varepsilon_s) = 0, s \neq t \\ E(x_s \varepsilon_t) = 0, \forall s < t \end{cases} \quad (2.1)$$

式中,  $\nabla^d = (1-B)^d$ ;  $\Phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$  为平稳可逆 ARMA( $p, q$ )模型的自回归系数多项式;  $\Theta(B) = 1 - \theta_1 B - \dots - \theta_q B^q$  为平稳可逆 ARMA( $p, q$ )模型的移动平均系数多项式。

### 2.2. 季节性差分自回归移动平均模型原理(SARIMA)

SARIMA 模型来源于自差分回归移动平均模型(ARIMA), 又称为季节乘积模型。在短期相关性和季节性影响的乘积关系中, 拟合模型实际上为 ARMA( $p, q$ )和 ARMA( $P, Q$ )<sub>S</sub> 相乘的结果。结合  $d$  阶趋势差分 and 以周期  $S$  为步长的  $D$  阶季节差分进行建模, 乘法模型的构造如下列公式所示[17]:

$$\nabla^d \nabla_S^D x_t = \frac{\Theta(B)\Theta_S(B)}{\Phi(B)\Phi_S(B)} \varepsilon_t \quad (2.2)$$

其中  $\Theta(B) = 1 - \theta_1 B - \dots - \theta_q B^q$ ,  $\Phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$ ,  $\Theta_S(B) = 1 - \theta_1 B^S - \dots - \theta_Q B^{QS}$ ,  $\Phi_S(B) = 1 - \phi_1 B^S - \dots - \phi_P B^{PS}$ , 该乘法模型简记为 ARIMA( $p, d, q$ ) $\times$ ( $P, D, Q$ )<sub>S</sub>。

## 3. 集合经验模态分解理论(EEMD)

### 3.1. 经验模态分解基本原理

经验模态分解(EMD), 是美国华裔科学家黄诺登博士于 1998 年提出的一种新的自适应信号时频域处理技术[18]。它根据数据本身的时间特性对信号进行分解为有限个本征模函数(IMF), 每个 IMF 分量分别为源信号的不同时间尺度的局部特征信号。

为了从原始序列中分解出 IMF, EMD 的分解过程如下:

先将数据序列分段筛选出极大值和极小值点, 上下极值点的包络线  $e_{\max(t)}$  和  $e_{\min(t)}$ , 用三次样条曲线拟合出来, 并计算上下包络线的平均值  $m(t)$ 。设原始数据为  $x(t)$ , 在  $x(t)$  中减去  $m(t)$  得到:

$$h(t) = x(t) - m(t) \quad (3.1)$$

再根据预设判据判断  $h(t)$  是否为 IMF, 重复以上过程直到  $h(t)$  满足判据, 则  $h(t)$  就是需要提取的 IMF:  $C_k(t)$ 。每得到一次 IMF 分量后从原信号中扣除, 直到剩余部分  $r_n(t)$  为单调的序列或者常数序列。这样原始序列可分解为:

$$x(t) = \sum_{i=1}^N C_i(t) + r_n(t). \quad (3.2)$$

### 3.2. 集合经验模态分解基本原理

集合经验模态分解(EEMD), 是为了解决 EMD 方法存在模态混叠等不足而提出的一种叠加高斯白噪

声的多次经验模态分解。利用高斯白噪声具有频率均匀分布的统计特性,通过每次加入同等幅值的不同白噪声来改变信号的极值点特性,再对多次经验模态分解得到的相应 IMF 进行总体平均来抵消加入的白噪声,从而有效抑制模态混叠的产生。

EEMD 分解的具体步骤[19][20]如下:

第一步,设定总体平均次数(集合数)  $m$ ;

第二步,将在原始序列  $x(t)$  中加入高斯白噪声序列  $\varepsilon_i(t), i=1,2,\dots,m$ 。

$$x_i(t) = x(t) + \varepsilon_i(t). \quad (3.3)$$

第三步,对含噪序列  $x_i(t)$  分别进行 EMD 分解,分解各自 IMF 分量  $C_{ij}(t), j=1,2,\dots,J$  和残差分量  $r_i(t)$ 。

$$x_i(t) = \sum_{j=1}^J C_{ij}(t) + r_i(t). \quad (3.4)$$

第四步,对于所得到对应的 IMF 分量求均值。从而得到原序列的第  $k$  个 IMF 分量  $c_k(t)$  和剩余分量  $r(t)$ :

$$c_k(t) = \frac{1}{m} \sum_{i=1}^m C_{ik}(t), r(t) = \frac{1}{m} \sum_{i=1}^m r_i(t) \quad (3.5)$$

第五步,原序列可分解为

$$x(t) = \sum_{k=1}^m c_k(t) + r(t) \quad (3.6)$$

#### 4. 奇异谱分析理论(SSA)

奇异谱分析(SSA) [21][22]是 1978 年 Colebrook 提出的一种用于非线性时序资料的新方法。建立了观察到的时间资料的轨迹矩阵,并将轨迹矩阵分解、重建,从表示原始时刻的序列中分离出各种分量信号,例如长期趋势、周期、噪声信号等,从所提取的信号中进一步分析和预测时间序列结构。

第一步,嵌入。将一维时间序列数据  $y = (y_1, y_2, \dots, y_N)$  转化为其轨迹矩阵  $X$ :

$$X = (x_{ij})_{i,j=1}^{L,K} = \begin{pmatrix} y_1 & y_2 & \cdots & y_K \\ y_2 & y_3 & \cdots & y_{K+1} \\ \vdots & \vdots & \ddots & \vdots \\ y_L & y_{L+1} & \cdots & y_N \end{pmatrix} \quad (4.1)$$

其中,  $L$  为选取的窗口长度,  $1 < L < N$ ,  $K = N - L + 1$ 。

第二步,奇异值分解。令  $S = XX^T$ , 对  $S$  进行奇异值分解后得到的其  $L$  个特征值  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_L \geq 0$  和其所有的正交特征向量  $U_1, U_2, \dots, U_L$ , 令  $d = \max\{i, \lambda_i > 0\}$ , 记  $V_i = \frac{X^T U_i}{\sqrt{\lambda_i}}, i=1,2,\dots,d$ 。则矩阵  $X$  的

奇异值分解可以写成:

$$X = X_1 + X_2 + \dots + X_d \quad (4.2)$$

其中,  $X_i = \sqrt{\lambda_i} U_i V_i^T$ ,  $\sqrt{\lambda_i}$  是矩阵  $X$  的奇异值,  $U_i$  是左特征向量,  $V_i$  是右特征向量,  $(\sqrt{\lambda_i}, U_i, V_i)$  称为矩阵  $X$  的第  $i$  个三重特征向量。

第三步,分组。以不同的提取成分作为依据,将  $X_i$  分为  $m$  个不同的组  $I_1, I_2, \dots, I_m$ , 并将每组的矩阵相加。如第  $I_i$  组包含的子集为  $I_i = \{i_1, \dots, i_p\}$ , 则

$$X_{I_i} = X_{i_1} + X_{i_2} + \dots + X_{i_p} \quad (4.3)$$

而  $X$  相应分解为

$$X = X_{I_1} + X_{I_2} + \dots + X_{I_m} \tag{4.4}$$

其中,  $X_{I_j}$  的贡献率可表示为  $\sum_{i \in I_j} \lambda_i / \sum_{1 \leq j \leq m, i \in I_j} \lambda_i$ 。

第四步, 重构。设  $Y$  为  $L \times K$  维矩阵, 矩阵元为  $y_{ij}$ ,  $1 \leq i \leq L, 1 \leq j \leq K$ 。定义  $L^* = \min(L, K)$ ,  $K^* = \max(L, K)$ ,  $N = L + K - 1$ ,  $y_{ij}^* = \begin{cases} y_{ij}, & \text{若 } L < K \\ y_{ji}, & \text{若 } L \geq K \end{cases}$ , 则重构序列  $G = (g_0, g_1, \dots, g_{N-1})$  可通过下列式子计算获得:

$$g_k = \begin{cases} \frac{1}{k+1} \sum_{m=1}^{k+1} y_{m, k-m+2}^* & \text{若 } 0 \leq k \leq L^* - 1 \\ \frac{1}{L^*} \sum_{m=1}^{L^*} y_{m, k-m+2}^* & \text{若 } L^* \leq k \leq K^* - 1 \\ \frac{1}{N - K^*} \sum_{m=k-K^*+2}^{N-K^*+1} y_{m, k-m+2}^* & \text{若 } K^* \leq k \leq N \end{cases} \tag{4.5}$$

式(4.5)本质上是对矩阵  $X_{ij}$  在对角线方向上求出  $i + j = k + 2$  各单元的平均, 并求得相应的  $g_k$  值, 从而获得  $X_{ij}$  的重组序列  $G$ 。

### 5. 小波变换理论(WT)

小波变换(Wavelet Transform, 缩写为 WT), 又称小波分析[23] [24] [25], 是指信号通过有限长或者快速衰减的“母小波”的振荡波形来表出, 它不仅传承和延续了短时傅里叶变换的局部化思想, 也克服和改进短时傅里叶变换的窗口大小不随频率变化等缺点。小波变化旨在对信号进行多尺度细化分析, 信号中的信息有效地通过伸缩和平移等运算的方式, 以及时间和领域的局域变换被提取出来。小波分析主要包括分解、去噪和重建三个步骤。

小波分解是将某一小波基函数  $\psi$  做位移  $b$  后, 再在不同尺度  $a$  下, 与尚未分析的信号  $y$  做内积。上述过程可逆, 其逆过程称作小波重构。小波分解和小波重构的表达式分别如下:

$$W_y(a, b) = \langle y(t), \Psi_{a,b} \rangle = |a|^{-1/2} \int_{-\infty}^{+\infty} \overline{\Psi\left(\frac{t-b}{a}\right)} y(t) dt \tag{5.1}$$

$$y(t) = K \iint_{-\infty}^{+\infty} W_y(a, b) |a|^{-1/2} \Psi\left(\frac{t-b}{a}\right) \frac{da db}{a^2} \tag{5.2}$$

其中,  $W_y(a, b)$  为小波分解系数,  $a$  为伸缩因子,  $b$  为平移因子,  $\Psi_{a,b}$  为小波基函数,  $\langle, \rangle$  为内积,  $y(t)$  为数据,  $\bar{\Psi}$  是  $\Psi$  的共轭。

### 6. BP 神经网络理论

反向传播神经网络(Back Propagation Neural Network, 简称 BP 神经网络)运用输入数据的训练让网络存在联想记忆, 从而使得联想记忆发挥其预测能力。模型的主要结构分为输入层, 隐含层, 输出层。BP 神经网络的结构如图 1 所示。信号通过输入层、隐含层向输出层的正向传播, 误差的逆向传播, 以及对预测的错误部分的误差调节, 使得预测结果与预期输出不断逼近[26]。

各个层次的数学关系如下所示:

对于隐含层有:

$$y_j = f(\text{net}_j), \text{net}_j = \sum_{i=1}^n v_{ij} x_i, j = 1, 2, \dots, m \tag{6.1}$$

对于输出层有:

$$o_k = f(\text{net}_k), \text{net}_k = \sum_{i=0}^m w_{ik} y_i, j = 1, 2, \dots, l \quad (6.2)$$

一般情形下激活函数  $f(x)$  取  $\frac{1}{1+e^{-x}}$ , 也可以根据需要取  $\frac{1-e^{-x}}{1+e^{-x}}$ 。

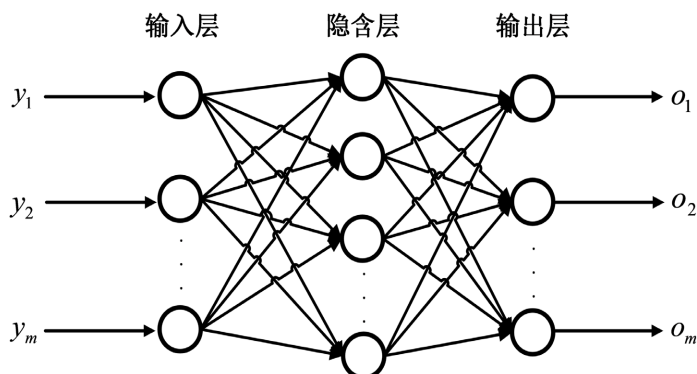


Figure 1. Structure of BP neural network

图 1. BP 神经网络的结构

## 7. 实证分析——全国铁路客运量预测分析

### 7.1. 数据描述

本文数据来自于国家统计局发布的月度数据报告中的铁路客运量当期值(万人)。本文选取 1984 年 1 月至 2022 年 1 月, 共 457 条铁路月客运量数据进行建模和检验。图 2 中可以看出, 铁路客运量当期量总体呈上升后近两年略有下降的趋势。由于新型冠状病毒疫情的爆发, 2020 年 2 月铁路客运量达到近十五年的最低值, 铁路客运量仅有约 3723 万人。而后至 2022 年 1 月, 铁路客运量的波动较大。

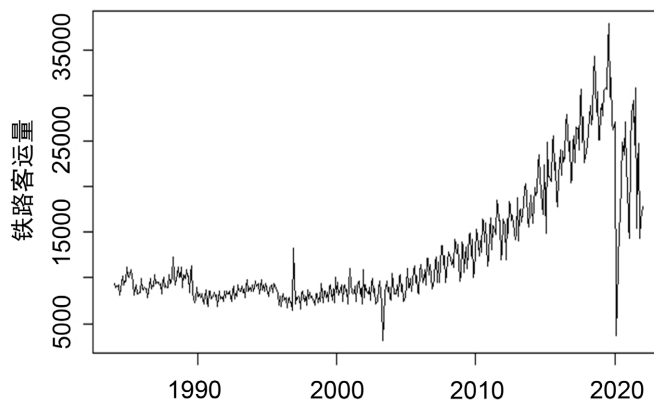


Figure 2. Time series of raw data

图 2. 原始数据时序图

### 7.2. 原始序列 SARIMA 预测结果

本文使用 R 语言程序对原始序列构建 SARIMA 模型对铁路月客运量进行预测。一阶差分后通过平稳性和白噪声检验。首先确定短期相关模型, 利用模型 ARMA(2,4) 作为原始序列差分后得到的短时自相关信息模型。考虑到季节自相关特性的情况下, 使用 ARMA(0,2)<sub>12</sub> 模型拟合差分序列的季节自相关信息。因此, 得到模型为 ARIMA(2,1,4)×(0,0,2)<sub>12</sub>。对拟合模型进行检验, 结果显示该模型通过残差白噪声检

验，因此拟合函数可以写为：

$$\nabla x_t \nabla^{12} = \frac{1 + 0.4000B + 0.6108B^2 - 0.2056B^3}{1 - 0.4628B^2} (1 - 0.2849B^{12} - 0.3373B^{24}) \varepsilon_t.$$

将 SARIMA 模型拟合得到的结果与铁路月客运量的实际值进行对比，如图 3 所示。

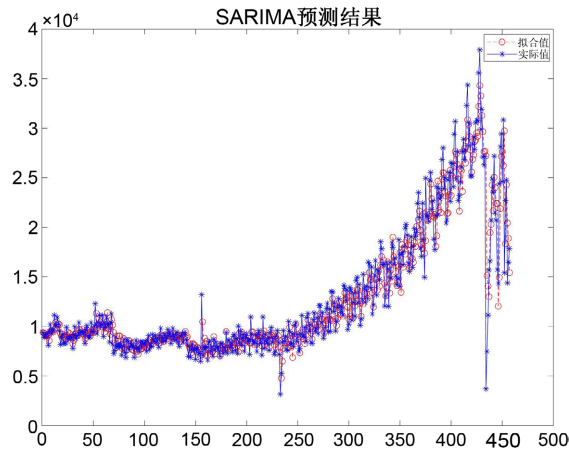


Figure 3. Comparison of SARIMA predictions  
图 3. SARIMA 预测结果对比图

### 7.3. 原始序列 BP 预测结果

初始化网络、训练和仿真是 BP 神经网络[25]在进行建模时的三个基本步骤。选取默认的函数线性函数('tansig')，根据经验公式指定隐含层神经元个数为 10。同时本文将最小均方误差(MSE)选作模型的评价指标，设置期望误差最小值为 0.0001，最大训练步长为 1000，学习率为 0.05。其余设置则采用 MATLAB 神经网络工具箱中的默认设置。训练函数选择 Levenberg-Marquardt 算法。同时设定 90% 的训练样本数，5% 的验证样本数和 5% 的测试样本数，即 411 个训练样本、23 个验证样本和 23 个预测样本。

将原始序列未经任何分解后直接通过上述设定的 BP 神经网络模型进行训练和仿真，得到的预测结果与实际值的对比如图 4 所示。

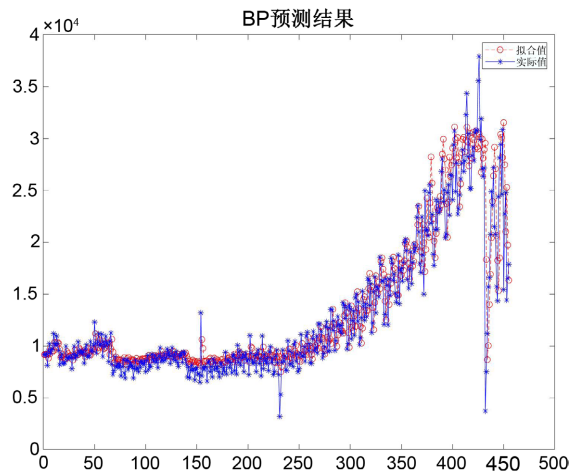


Figure 4. Comparison of SARIMA predictions  
图 4. SARIMA 预测结果对比图

## 7.4. EEMD 分解预测模型

第一步,运用 EEMD 分解法对原始数据进行分解。Nstd 是设置高斯白噪声的标准差,本文取 0.2。NE 是添加噪声的次数,本文设置为 100。将加入白噪声序列的铁路月客运量的原始序列进行分解,得到 7 个 IMF 分量和 1 个剩余量,根据频率高到低的顺序将分解得到的分量排序,原始时间序列的波动特征成份由各自的波动代表,剩余量为无法分解的部分,如图 5 所示。

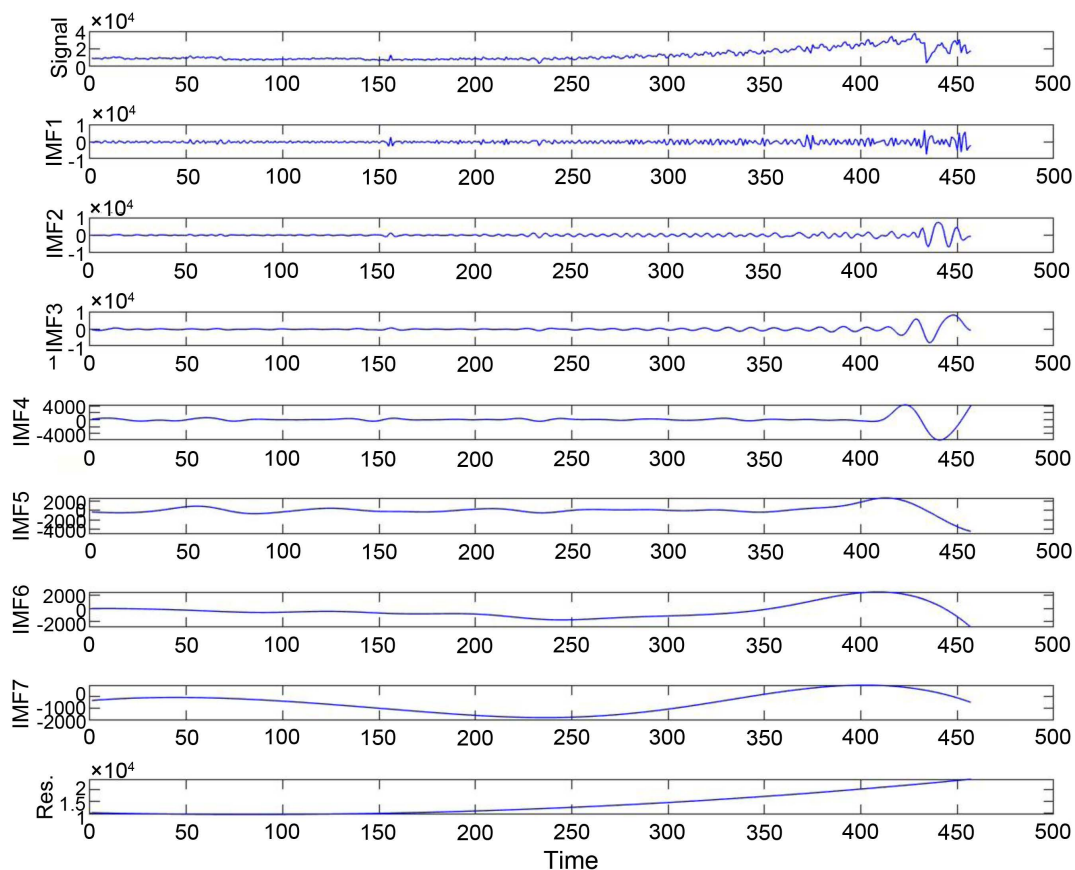


Figure 5. Time series of EEMD decomposed IMF components and residual components

图 5. EEMD 分解后的 IMF 分量及剩余分量时序图

第二步,对 EEMD 分解后的 7 个 IMF 分量和残差项分别进行 SARIMA 时间序列预测。SARIMA 模型定阶全部由 R 语言程序中 forecast 包中的 auto.arima()函数所选取的 AIC 最小的模型。得到 7 个 IMF 和残差项的模型参数如表 1 所示。

Table 1. Model order corresponding to each component

表 1. 各分量对应的模型阶数

	ARIMA			SEASONAL		
	$p$	$d$	$q$	$p$	$d$	$q$
IMF1	2	0	1	0	0	1
IMF2	4	0	2	0	0	1
IMF3	4	0	4	0	0	2



Continued

IMF4	1	0	0	1	0	0
IMF5	0	0	0	1	0	0
IMF6	0	2	0	0	0	0
IMF7	0	2	0	0	0	0
RES	0	2	5	0	0	2

第三步, 将 EEMD 分解后得到的 7 个 IMF 分量和残差项分别通过 BP 神经网络模型进行预测[19]。隐含层神经元个数为 10, 网络参数设定与训练样本比例设定与前文相同, 得到 BP 预测结果。

通过对 EEMD 分解后的 7 个 IMF 分量和残差项分别通过 SARIMA 模型和 BP 神经网络预测模型进行预测后, 以 MAPE 为指标, 得到预测效果如表 2 所示。从表 2 中可以看出, 对于 IMF1、IMF2、IMF4、IMF5、IM6、IM7, BP 神经网络的预测效果较好。对于 IMF3, SARIMA 模型具有更高的精度。而残差项两者的预测误差都很小。

**Table 2.** Comparison of forecast results (MAPE)

**表 2.** 预测结果对比(MAPE)

	IMF1	IMF2	IMF3	IMF4	IMF5	IMF6	IMF7	RES
EEMD-SARIMA	131.29%	37.56%	4.01%	32.80%	139.59%	2.74%	1.39%	0.00%
EEMD-BP	48.51%	16.10%	15.33%	2.37%	0.09%	0.05%	0.11%	0.00%

最后, 进行 EEMD-SARIMA-BP 组合预测。根据表 2, 对 IMF3 和残差项选取 SARIMA 模型进行预测, 对其他分量选取 BP 神经网络方法进行预测, 再将得到的预测的分量值叠加后进行组合预测。得到组合预测模型的 MAPE 为 3.96%, 提高了预测精度。

取 2020 年 3 月至 2022 年 1 月共 23 个月的铁路月客运量的值作为预测组, 将原始序列直接采用 SARIMA 模型预测、直接采用 BP 神经网络模型预测、EEMD 分解后分别采用 SARIMA 模型预测、EEMD 分解后分别采用 BP 模型预测和 EEMD 分解后 SARIMA 和 BP 模型组合预测得到结果的 MAE、MAPE、RMSE 进行对比, 得到表 3。

**Table 3.** Precision comparison of EEMD decomposition prediction model

**表 3.** EEMD 分解预测模型精度对比

	BP	SARIMA	EEMD-SARIMA	EEMD-BP	EEMD-SARIMA-BP
MAE	5303.2003	3907.7889	2240.0609	935.9701	811.9572
MAPE	31.93%	23.51%	13.39%	4.35%	3.96%
RMSE	6515.2181	5044.6015	2924.4861	1202.0687	1120.4639

## 7.5. SSA 分解预测模型

第一步, 采用 Matlab 实现奇异谱分解[22]。首先, 将数据标准化处理。再设定窗口长度 L (嵌入维数), 一般设定的长度小于或者等于序列长度的一半, 将原始序列延时排列成矩阵形式。本文的铁路月客运量数据为 457 条, 将 L 确定为 200, 按非增的次序求出 200 个特征值。第一和第二个特征值被采用(贡献率为 88.55%), 见图 6。并用公式(4.5)计算其对应的重组序列, 由此获得到低频段结果。剩余的特征值来计算延时矩阵, 同用公式(4.5)求出对应的重构序列, 从而得到高频部分。原始序列和分解成份的时序图如

图 7 所示,可以看出铁路月客运量整体先缓慢上升后加速上升后开始逐步下降的趋势。而高频率序列整体的波动情况较为稳定,随机波动性较强的部分可视为噪声成分,但由于内在蕴含着一定波动变化,不完全归于噪声,其中 2020 年左右的波动较其它时间更加剧烈,其代表了铁路月客运量序列中的随机性成分。

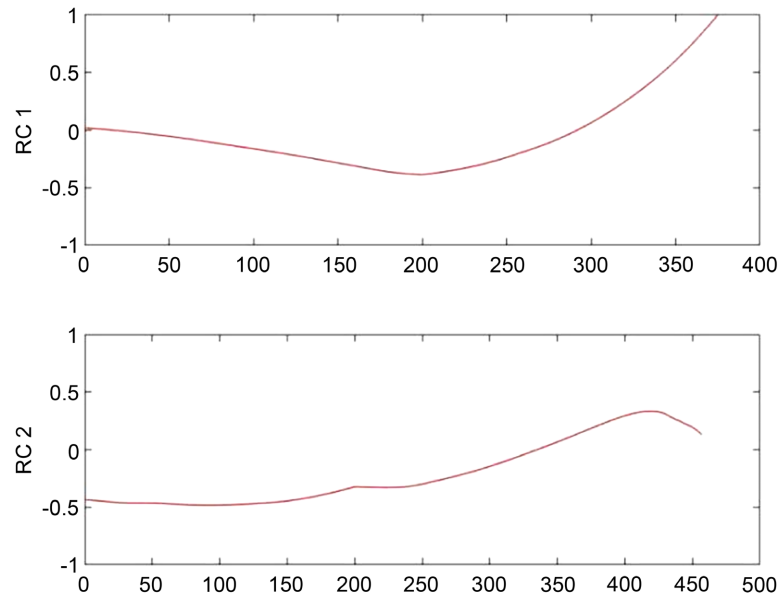
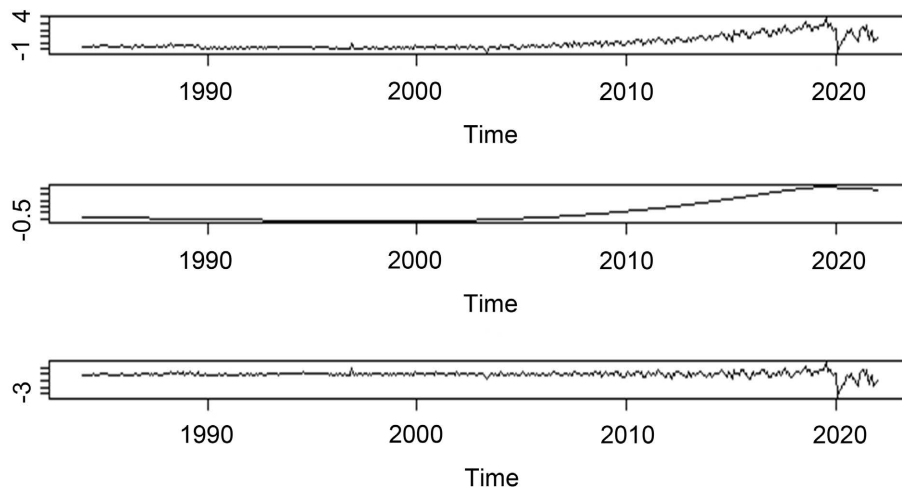


Figure 6. Low frequency sequence

图 6. 低频率序列



注: 从上至下为原始序列、重构的低频率序列、重构的高频率序列,数据均标准化处理。

Figure 7. Original series and series diagram reconstructed after decomposition

图 7. 原始序列与分解后重构的时序图

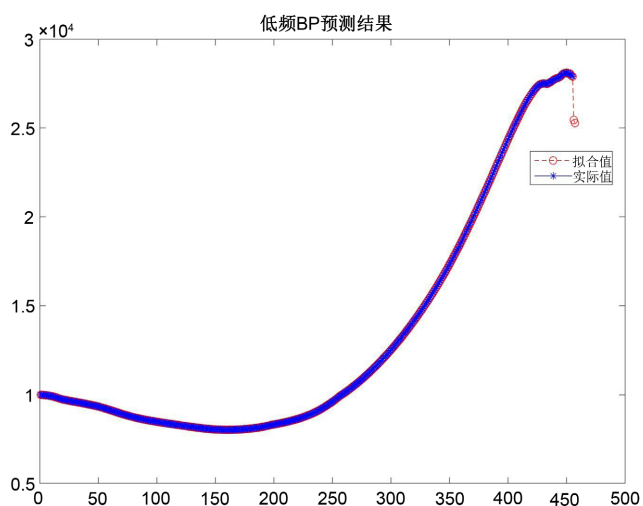
第二步,对分解得到的低频和高频分量分别进行 SARIMA 模型建模,得到高低频分量的模型参数如表 4 所示。取低频项和高频项分别预测的 2020 年 3 月至 2022 年 1 月共 23 个月的铁路客运量的值进行对比,得到平均绝对百分误差 MAPE 的值分别为 11.11% 和 95.02%。

第三步,分别对低频部分和高频部分进行 BP 神经网络时间序列预测。模型参数设置与前文相同。图 8 为低频部分的预测结果,图 9 为高频部分的预测结果。同样选取 2020 年 3 月至 2022 年 1 月的客

运量实际值和预测值进行平均相对误差的比较, 得到低频分量的 MAPE 为 0.04%, 而高频分量的 MAPE 为 86.05%。

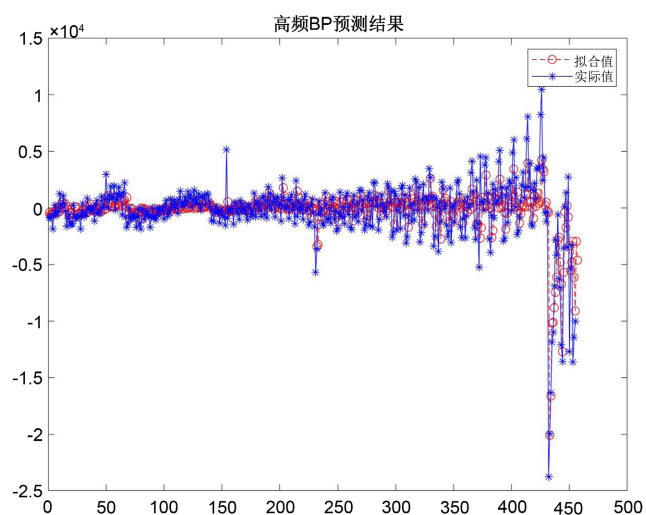
**Table 4.** Model order determination parameters of high and low frequency components  
**表 4.** 高低频分量的模型定阶参数

	ARIMA			SEASONAL		
	$p$	$d$	$q$	$p$	$d$	$q$
低频分量	4	2	1	1	0	2
高频分量	3	0	1	1	0	0



**Figure 8.** BP prediction of low frequency series

**图 8.** 低频部分的 BP 预测结果



**Figure 9.** BP prediction of high frequency series

**图 9.** 高频部分的 BP 预测结果

相对于奇异谱分解的 SARIMA 模型预测, 分解后采用 BP 神经网络模型预测得到的高低频分量的误差都较小。因此分解后采用 BP 模型进行预测可以有效地减小误差, 不需要进行组合预测。

取 2020 年 3 月至 2022 年 1 月共 23 个月的铁路月客运量的值作为预测组, 将原始序列直接采用 SARIMA 模型预测、直接采用 BP 神经网络模型预测、SSA 分解后分别采用 SARIMA 模型预测、SSA 分解后分别采用 BP 模型预测得到结果的 MAE、MAPE、RMSE 进行对比, 结果如表 5 所示。

**Table 5.** Precision comparison of SSA decomposition prediction model

**表 5.** SSA 分解预测模型精度对比

	BP	SARIMA	SSA-SARIMA	SSA-BP
MAE	5303.2003	3907.7889	3513.6369	2908.2811
MAPE	31.93%	23.51%	23.11%	15.65%
RMSE	6515.2181	5044.6015	4853.0374	3960.8071

## 7.6. 小波分解预测模型

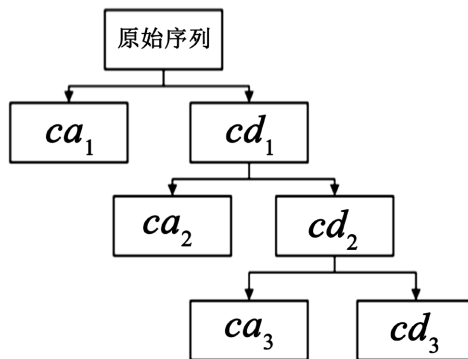
第一步, 进行序列的小波分解与重构。由于没有标准划一的方法来选定小波基函数, 因此本文经过尝试后选择具有正交性和紧支撑性[23]的 db5 小波(即 Daubechies 小波)作为小波函数。若小波分解的层数太少, 会影响到降噪后的平滑度和平稳性, 而分解的层数太多, 则导致了在分解时的运算错误较大[24], 所以在综合考虑后, 本文采用了三层小波分解法。图 10 显示了分解的结构。用 Matlab 软件进行分解和重组, 其结果如图 11 所示。其中  $a_3$  为低频序列  $ca_3$  重构后的分量, 基本保持了原始序列的形状, 反映了铁路月客运量的长期变化趋势。 $d_1$ 、 $d_2$  和  $d_3$  为重构后的高频序列分量, 其隐藏了原始序列的细节特征, 反映了铁路客运量的波动特征和随机特征。

第二步, 对小波分解及重构后得到的低频序列  $a_3$ 、高频序列  $d_1$ 、高频序列  $d_2$ 、高频序列  $d_3$  分别进行 SARIMA 模型建模, 得到 4 个分量的模型参数如表 6 所示。

第三步, 将小波分解重构后得到的高低频数据 BP 神经网络训练和仿真。图 12 为经过小波分解重构后的序列进行分别进行 BP 神经网络预测叠加后的拟合图, 可以看出拟合效果较好。

通过对小波分解重构后的 4 个序列分别通过 SARIMA 模型和 BP 模型进行预测后, 同样选取 2020 年 3 月至 2022 年 1 月的客运量实际值和预测值进行平均相对误差的比较, 得到预测效果如表 7 所示。从表 7 中可以看出, 对于低频序列  $a_1$ , BP 模型和 SARIMA 模型预测精度差别不大, 都有较小误差。而对于  $d_1$ 、 $d_2$ , SARIMA 模型具有更高的精度。对于  $d_3$ , BP 模型在预测准确性上有优势。

最后进行 WT-SARIMA-BP 组合预测。对高频序列  $d_1$  和  $d_2$  选取 SARIMA 模型进行预测, 对低频序列  $a_1$  和高频序列  $d_3$  选取 BP 神经网络方法进行预测, 再将得到的预测的分量值叠加后进行组合预测。



**Figure 10.** Structure of three-layer wavelet decomposition

**图 10.** 三层小波分解结构图

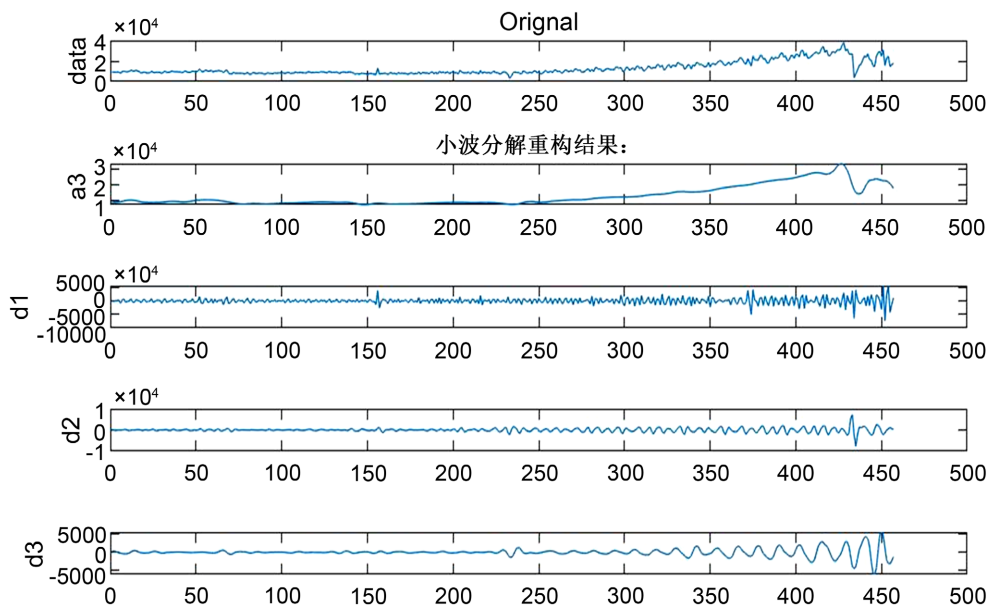


Figure 11. Results of wavelet decomposition and reconstruction

图 11. 小波分解与重构后的结果

Table 6. Model order determination parameters of high and low frequency components

表 6. 高低频分量的模型定阶参数

	ARIMA			SEASONAL		
	$p$	$d$	$q$	$p$	$d$	$q$
a3	2	1	0	0	0	2
d1	1	0	0	1	0	1
d2	5	0	0	2	0	0
d3	5	0	0	0	1	2

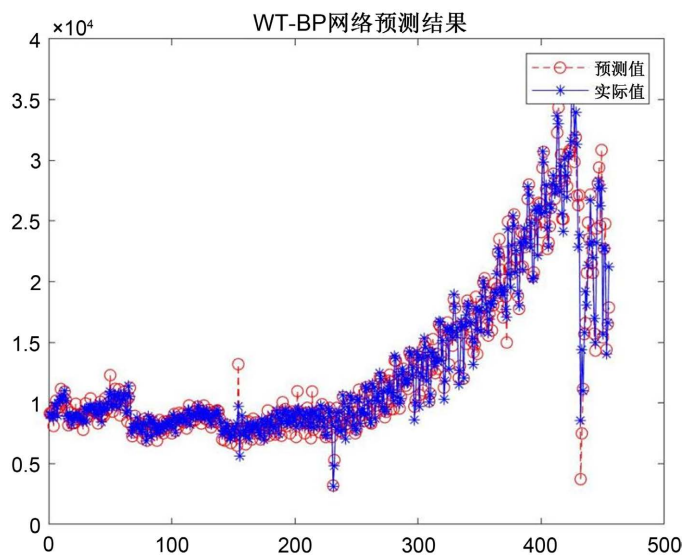


Figure 12. WT-BP forecast result

图 12. WT-BP 预测结果图

**Table 7.** Comparison of wavelet decomposition prediction results**表 7.** 小波分解预测结果对比

	a1	d1	d2	d3
WT-SARIMA	1.61%	81.35%	40.33%	27.76%
WT-BP	1.65%	683.51%	62.57%	18.38%

取 2020 年 3 月至 2022 年 1 月共 23 个月的铁路月客运量的值作为预测组, 将原始序列直接采用 SARIMA 模型预测、直接采用 BP 神经网络模型预测、WT 分解后分别采用 SARIMA 模型预测、WT 分解后分别采用 BP 模型和 WT 分解后采用组合的 SARIMA 和 BP 模型预测得到结果的 MAE、MAPE、RMSE 进行对比, 得到表 8。组合预测模型得到预测值的 MAPE 为 9.72%, 相较于 WT-SARIMA 模型的 12.83% 与 WT-BP 模型的 9.80% 有所下降。

**Table 8.** Precision comparison of wavelet decomposition prediction models**表 8.** 小波分解预测模型精度对比

	BP	SARIMA	WT-SARIMA	WT-BP	WT-SARIMA-BP
MAE	5303.2003	3907.7889	2425.5059	1534.7437	1720.2375
MAPE	31.93%	23.51%	12.83%	9.80%	9.72%
RMSE	6515.2181	5044.6015	3120.4066	2199.2559	2284.4745

## 7.7. 预测结果对比

### 7.7.1. 分解方法预测精度对比

对于 BP 神经网络模型, 通过对比 EEMD、SSA 和小波分解的预测精度指标, 得到表 9。从表 9 中可以看出, 三种分解方法都可以提高 BP 神经网络预测的精度。对于 BP 预测模型, 综合三个指标的数据分析, 精度从高到低的排序为: EEMD-BP 模型、WT-BP 模型、SSA-BP 模型、BP 模型。

对于 SARIMA 模型, 通过对比三种分解方法的预测精度指标, 得到表 10。从表 10 可以看出, 经过采用三种分解方法均可以提高 SARIMA 模型的预测精度。对于 SARIMA 预测模型, 综合三个指标的数据分析, 精度从高到低的排序为: WT-SARIMA 模型、EEMD-SARIMA 模型、SSA-SARIMA 模型、SARIMA 模型。

**Table 9.** Accuracy of BP prediction models with different decomposition methods**表 9.** 不同分解方法用法 BP 模型预测的精度指标

	BP	EEMD-BP	SSA-BP	WT-BP
MAE	5303.2003	935.9701	2908.2811	1534.7437
MAPE	31.93%	4.35%	15.65%	9.80%
RMSE	6515.2181	1202.0687	3960.8071	2199.2559

**Table 10.** Accuracy of SARIMA prediction models with different decomposition methods**表 10.** 不同分解方法用 SARIMA 模型预测的精度指标

	SARIMA	EEMD-SARIMA	SSA-SARIMA	WT-SARIMA
MAE	3907.7889	2240.0609	3513.6369	2425.5059
MAPE	23.51%	13.39%	23.11%	12.83%
RMSE	5044.6015	2924.4861	4853.0374	3120.4066

### 7.7.2. 所有预测方法精度对比

对上述共 10 种不同的预测模型的精度指标从小到大进行大致排名, 得到表 11。从表 11 可以总结出: 对于铁路月客运量数据, 采用 EEMD、SSA 和 WT 三种分解方法都可以提高预测的精度。采用 EEMD 分解方法的预测效果普遍优于采用 WT 分解, 而 WT 分解法的预测效果又优于 SSA 分解法。分解后采用 BP 模型进行预测的误差普遍小于分解后采用 SARIMA 模型预测得到的误差。其中, 预测效果最好的是采用 EEMD-SARIMA-BP 组合模型进行预测, 而预测效果最差的是直接使用 BP 模型进行预测。

Table 11. Accuracy of all prediction methods

表 11. 所有预测方法的精度指标

	MAE	MAPE	MRSE
EEMD-SARIMA-BP	811.9572	3.96%	1120.4639
EEMD-BP	935.9701	4.35%	1202.0687
WT-SARIMA-BP	1720.2375	9.72%	2284.4745
WT-BP	1534.7437	9.80%	2199.2559
WT-SARIMA	2425.5059	12.83%	3120.4066
EEMD-SARIMA	2240.0609	13.39%	2924.4861
SSA-BP	2908.2811	15.65%	3960.8071
SSA-SARIMA	3513.6369	23.11%	4853.0374
SARIMA	3907.7889	23.51%	5044.6015
BP	5303.2003	31.93%	6515.2181

## 8. 结论

铁路客运量具有较强的季节性, 也具有趋势性, 同时在内外部因素影响下还具有随机性和非线性等特点。因此, 本文引入的具有季节效应的自回归移动平均模型和具有强非线性模拟、自学习能力的神经网络方法可对客运量进行有效的预测。基于改进预测精度的想法, 本文引入集合经验模态分解、奇异谱分解和小波分解三种方法先对序列数据进行处理后再进行不同预测模型的建模。对 1984 年 1 月至 2022 年 1 月的全国铁路月客运量序列数据进行训练和预测, 得到了以下结论:

1) 通过集合经验模态分解、奇异谱分解和小波分解三种分解方法处理过的数据, 得到不同频率的分量再建立预测模型进行预测均可以有效地提高序列的预测精度。这说明分解数据可以为预测提取更有效的信息, 使得预测模型的建模更具有针对性。

2) 在集合经验模态分解方法中, EEMD-SARIMA-BP 组合模型的预测精度是本文中提到的十种预测模型中最好的。而 EEMD-BP 模型的精度略低于组合模型, 这是因为组合模型中部分分量采用 SARIMA 可以得到更佳的预测效果, 因此叠加后的组合模型可以提高整体预测的精确度。若不进行组合预测, 将所有的 IMF 分量和残余项均进行 BP 模型预测, 得到的预测结果的误差则是小于将所有分量全部进行 SARIMA 模型预测的误差。

3) WT-SARIMA-BP 组合模型是所有小波分解预测方法中精度最高的, 其次是 WT-BP 模型, 预测效果最差的是 WT-SARIMA 模型。这说明, BP 神经网络模型进行预测较 SARIMA 方法更适用于小波分解与重构后的铁路月客运量序列, 可以提高预测的精确性。虽然 WT-SARIMA 模型的精度略低于上述两种方法, 但还是高出原始序列直接采用 SARIMA 模型和 BP 模型预测的精度很多。

4) 虽然预测的精度低于上述两种分解方法的预测模型, 但 SSA-BP 模型和 SSA-SARIMA 的预测精度也是高于直接用预测模型进行预测的。这表明通过奇异谱分解后的序列数据预测值较直接使用预测模型得到的预测值更加接近真实值, 用该种分解方法进行分解后再预测同样可以提高预测的精度, 但是对于铁路客运量预测值提升精度的效果有限。

## 参考文献

- [1] 国家铁路局介绍铁路领域“十三五”发展成就[J]. 铁道技术监督, 2020, 48(11): 53.
- [2] 罗庆中, 李娜, 贾光智. 中国铁路发展战略研究[J]. 科技导报, 2020, 38(9): 26-31.
- [3] 马广文, 王西秩总. 交通大辞典[M]. 上海: 上海交通大学出版社, 2005.
- [4] 侯福均, 吴祈宗. BP 神经网络在铁路客运市场时间序列预测中的应用[J]. 运筹与管理, 2003(4): 73-75.
- [5] 郝军章, 崔玉杰, 韩江雪. 基于 SARIMA 模型在我国铁路客运量中的预测[J]. 数学的实践与认识, 2015, 45(18): 95-104.
- [6] 王雷, 金勇, 刘岩. 铁路客运量预测模型对比分析[J]. 山东交通学院学报, 2020, 28(3): 25-32+47.
- [7] 缪巧芬, 唐国强, 罗耀宁. 基于 X-13A-S 季节调整方法的铁路客运量预测分析[J]. 桂林理工大学学报, 2018, 38(3): 579-584.
- [8] 谢小山. 基于遗传算法和 BP 神经网络的铁路客运量预测研究[D]: [硕士学位论文]. 成都: 西南交通大学, 2010.
- [9] 王卓, 王艳辉, 贾利民, 李平. 改进的 BP 神经网络在铁路客运量时间序列预测中的应用[J]. 中国铁道科学, 2005(2): 130-134.
- [10] 汪健雄, 刘春煌, 单杏花, 朱建生. 基于双层次正交神经网络模型的铁路客运量预测[J]. 中国铁道科学, 2010, 31(3): 126-132.
- [11] 李万, 冯芬玲, 蒋琦玮. 改进粒子群算法优化 LSTM 神经网络的铁路客运量预测[J]. 铁道科学与工程学报, 2018, 15(12): 3274-3280.
- [12] 彭珍瑞, 孟建军, 祝磊, 蒋兆远. 基于支持向量机的铁路客运量预测[J]. 辽宁工程技术大学学报, 2007(2): 269-272.
- [13] 侯丽敏, 马国峰. 基于灰色线性回归组合模型铁路客运量预测[J]. 计算机仿真, 2011, 28(7): 1-3+30.
- [14] 贺晓霞, 鲍学英, 王起才, 董朝阳. 基于 GM-周期扩展组合模型的铁路客运量预测[J]. 铁道科学与工程学报, 2015, 12(3): 685-689.
- [15] 刘琳玥. 基于 PCA-BP 神经网络的铁路客运量预测模型研究[J]. 综合运输, 2016, 38(8): 43-47+73.
- [16] 王国贤, 范英兵, 王凤玲, 谢安琪. 时间序列和神经网络下我国铁路客运量的预测研究[J]. 黑河学院学报, 2021, 12(5): 182-185.
- [17] 易丹辉, 王燕. 应用时间序列分析[M]. 第五版. 北京: 中国人民大学出版社, 2019.
- [18] 黄诚惕. 希尔伯特-黄变换及其应用研究[D]: [硕士学位论文]. 成都: 西南交通大学, 2006.
- [19] 于群, 朴在林, 胡博. 基于 EEMD 和 BP 神经网络的短期光伏功率预测模型[J]. 电网与清洁能源, 2016, 32(7): 132-137.
- [20] 方雪清, 吴春胤, 俞守华, 张大斌, 欧阳庆. 基于 EEMD-LSTM 的农产品价格短期预测模型研究[J]. 中国管理科学, 2021, 29(11): 68-77. <https://doi.org/10.16381/j.cnki.issn1003-207x.2019.0765>
- [21] 程爽. 基于奇异谱分析的时间序列互相关分析[D]: [硕士学位论文]. 大连: 辽宁师范大学, 2018.
- [22] 周天清. 基于奇异谱分析的金融时间序列自适应分解预测研究[D]: [硕士学位论文]. 南昌: 华东交通大学, 2012.
- [23] 彭乃驰, 党婷. 基于小波分析的 BP-SARIMA 模型的 CPI 预测[J]. 统计与决策, 2018, 34(16): 22-25.
- [24] 章浙涛. 小波分析理论及其在变形监测中的应用研究[D]: [硕士学位论文]. 长沙: 中南大学, 2014.
- [25] 李腾. 基于小波分析和 BP 神经网络相结合的股票波动预测方法研究[D]: [硕士学位论文]. 天津: 天津大学, 2018.
- [26] 刘天, 姚梦雷, 黄继贵, 陈红缨, 黄淑琼, 杨雯雯, 蔡晶, 吴然. BP 神经网络在传染病时间序列预测中的应用及其 MATLAB 实现[J]. 预防医学情报杂志, 2019, 35(8): 812-816+821.