

# 基于EMD-LSTM模型的城市轨道交通进站客流量预测

郭以朋

山东交通学院轨道交通学院, 山东 济南

收稿日期: 2023年6月9日; 录用日期: 2023年10月31日; 发布日期: 2023年11月9日

## 摘要

城市轨道交通进站旅客的流量存在非线性、随机性、周期性等特点,导致长短时记忆神经网络(LSTM, Long Short-Term Memory)在有效处理上述数据特性时,其算法性能上对预置参数具有很大的时间依赖性,导致预测性能的下降。为了进一步提高模型的预测精度,研究中使用EMD经验模态分解,将客流时间信息的各个尺度分解为若干个固有模态分量(IMF),从而去除一些噪声影响;再使用LSTM神经网络对各固有模态分量进行学习和预测。预测结果显示,EMD-LSTM的组合模型相对于LSTM模型,其客流量的预测结果具有更高的精度,为城市轨道交通进站客流量的预测起到重要的参考作用。

## 关键词

短时客流数据, 经验模态分解, 长短时记忆神经网络, 进站客流预测

# Prediction of Passenger Flow of Urban Rail Transit into the Station Based on EMD-LSTM

Yipeng Guo

School of Rail Transportation, Shandong Transportation University, Jinan Shandong

Received: Jun. 9<sup>th</sup>, 2023; accepted: Oct. 31<sup>st</sup>, 2023; published: Nov. 9<sup>th</sup>, 2023

## Abstract

The flow of passengers entering urban rail transit stations exhibits characteristics such as nonlinearity, randomness, and periodicity. These traits pose challenges for the Long Short-Term Memory (LSTM) neural network, as its algorithm performance is highly dependent on preset parameters when effectively dealing with such data characteristics, leading to a decline in prediction performance. In order to enhance the accuracy of the model's predictions, the study employs Empiri-

cal Mode Decomposition (EMD) to decompose the time information of passenger flow into several Intrinsic Mode Functions (IMFs), thereby removing some noise interference. The LSTM neural network is then used to learn and predict the various IMFs. The results of the predictions demonstrate that the combined EMD-LSTM model, compared to the LSTM model alone, achieves higher accuracy in predicting passenger flow, thereby providing important references for predicting passenger flow in urban rail transit stations.

## Keywords

Short-Term Passenger Flow Data, Empirical Modal Decomposition, Long Short-Term Memory Neural Network, Inbound Passenger Flow Forecast

Copyright © 2023 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

地铁交通网络发展起来后,各线路之间的相互关联度会大大增强,客流的波动性和叠加性会使掌握网络客流规律更加困难。在城市轨道交通运行的过程中,经常会出现拥挤的现象,尤其是工作日的早晚高峰期,人流量更大,会更为拥堵。因此如何科学准确地预测进站客流量,并制定合理的应对措施,对提高城市轨道交通服务质量具有重要意义。如何利用历史数据对轨道交通进站客流量进行快速有效地预测,一直是城市轨道交通领域的研究热点问题。对于预测方法,有历史同期客流可参照且预测期较短的客流预测,方法有历史同比法和指数平滑法等。历史同比法简单,考虑了城市轨道交通网络结构的变化但是准确性相对较低,而指数平滑法稍微复杂,考虑了历史数据波动情况,能够应对历史客流中的突发性变化,预测期不长,准确性较高,但是不适合网络结构发生变化情况下的客流预测。与传统的预测方法相比,机器学习具有很强的非线性拟合能力以及处理大批量的数据并且找出内在的规律对未知的数据进行预测能力。比较典型的机器学习方法包括粒子群算法、模拟退火算法, BP 神经网络, SVM 支持向量机,然而以上的方法并不能很好地预测时间序列模型,常用循环神经网络(RNN)进行序列模型预测,但是由于 RNN 所使用的是反向传导算法,会导致激活函数累乘,从而很容易出现梯度消失或梯度爆炸,使用 LSTM 神经网络单元来替代 RNN 单元,该方法可以有效地解决 RNN 梯度爆炸与消失的问题。然而 LSTM 神经网络模型与其他神经网络模型相似,其中的某些参数,如时间窗大小、批号、隐层细胞数等,都取决于经验值和反复的实验验证结果,因此其效率较低下。隐层和隐层单元的数量直接决定着模型的拟合能力。时间窗大小和批处理大小直接影响着建模的训练过程和效果,参数不同训练出模型的预测性能有很大的差异。依靠经验设定的参数会导致预测结果具有较大的不稳定性,降低了预测精度。由于实际客流的复杂性,使用单一的神经网络方法进行预测可能是片面的。对于不同的交通条件,结合合适的模型可以整合各个模型的优势,从而获得更好的预测结果。根据上述研究模型的缺陷,本文提出一种短时客流预测的方法—构建 EMD-LSTM 组合预测模型,EMD 适用于处理非线性时间序列的问题,其将原始时间序列分解为 IMF 和残差,以达到降噪的目的,从而能够更加准确的预测城市轨道交通客流。

## 2. 相关工作

### 2.1. 经验模态分解法(EMD)

1998 年, Norden E. Huang 等人提出了一种信号处理模型——经验模态分解(EMD)。在该模型中,原

始时间序列数据通过分解得到的分量即本征模函数(Intrinsic Mode Functions, IMF)和残差。[1] [2]确定序列是否为 IMF 分量的主要前提条件有二个：第一，各 IMF 分量的极值点与过零点的数目必须相等，或者最多相差一个。第二，由序列的局部极大值的包络线与局部极小值包络线的平均值为零。[3] EMD 与传统分解算法的不同处在于处理非平稳性数据上具有明显的优势，在处理平稳性差的时间序列的时候其分解得到的 IMF 分量相比于原始数据，其波动特性更加平稳有规律，具有很好的鲁棒性，可以有效减少模态混叠对序列分解的干扰性。因此，经验模态分解法因其优势被应用于各大工程领域中，[4] [5] [6]并取得了极好的效果。

EMD 算法的主要计算过程如下：

- 1) 对于时间序列  $x(t)$ ，求得其中所有的局部极大值与极小值点；
- 2) 用三次样条插值拟合出数据的上下包络线，原始的数据被上下包络线所包络，如图 1 所示。并求上下包络线的平均值  $m(t)$ ，用原始数据  $x(t)$ 减去均值曲线  $m(t)$ 得到新的序列  $h(t)$ ，即：

$$h(t) = x(t) - m(t) \quad (1)$$

- 3) 判断  $h(t)$ 是否满足 IMF 的条件，如果  $h(t)$ 符合 IMF 分量条件，则认为  $h(t)$ 为该循环次数下分解得到的 IMF 分量，如果不是，则用  $h(t)$ 代替  $x(t)$ ，将  $h(t)$ 作为时间序列转入步骤(1)重新计算；
- 4) 每得到一次循环的  $h(t)$ ，从原信号  $x(t)$ 中将其去除，直到序列剩余部分  $r(t)$ 为单调；
- 5) 最后通过上述步骤，得到原始序列的分解结果如式

$$x(t) = \sum_{i=1}^n c_i(t) + r(t) \quad (2)$$

## 2.2. 长短期记忆神经网络(LSTM)

Hochreiter 和 Schmidhuber 研究提出了长短期记忆(Long Short-Term Memory, LSTM)神经网络单元来替代 RNN 单元，该方法可以有效地解决 RNN 梯度消失的问题。LSTM 通过长短期存储技术的“门”逻辑控制决定该数据是更新还是舍弃，让 LSTM 神经网络能够记得相对比较长期的信息，从而解决了梯度消失的问题。

基于在这些方面的优异表现，应用其对时间序列进行研究受到广泛关注。LSTM 的基本架构包括了一个记忆单元状态(Cell)和三个门，即遗忘门、输入门、输出门。[7]使用 Cell 和三个门结构能够解决传统循环神经网络的梯度消失问题，能有效反应历史信息对当前状态的影响。LSTM 的基本结构框架图如图 1 所示。

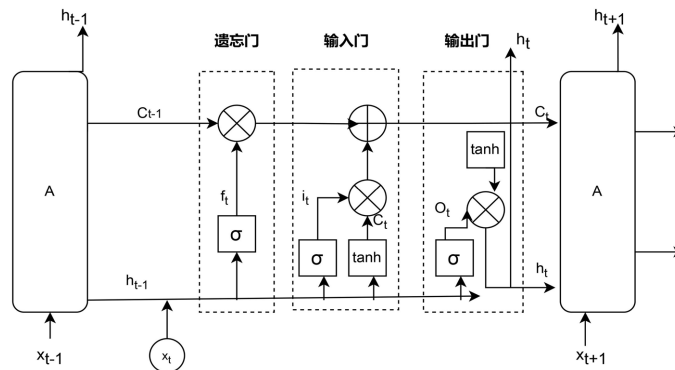


Figure 1. Schematic diagram of LSTM neural network structure  
图 1. LSTM 神经网络结构示意图

标准的 LSTM 架构, 整个网络由若干个同构单元(Cell)所组成, 包含的三种输入输出特征分别是:  $x_t$ 、 $h_{t-1}$ 、 $c_{t-1}$ , 输出分别为  $h_t$ 、 $c_t$ 。 $x_t$  代表  $t$  时刻的输入序列特征,  $h_{t-1}$  代表上一轮的状态量输出,  $c_{t-1}$  则代表上一轮全局一个信息的载体;  $h_t$  代表本轮的状态量输出,  $c_t$  则代表了本轮全局的一种信息载体。

[8] [9] LSTM 的每个单元可以用下列式子来表示:

$$\begin{aligned} f_t &= \sigma(W_f * [h_{t-1}, x_t] + b_f), \\ i_t &= \sigma(W_i * [h_{t-1}, x_t] + b_i), \\ c'_t &= \tanh(W_c * [h_{t-1}, x_t] + b_c), \\ c_t &= f_t * c_{t-1} + i_t * c'_t, \\ o_t &= \sigma(W_o * [h_{t-1}, x_t] + b_o), \\ h_t &= o_t * \tanh(c_t). \end{aligned} \quad (3)$$

以上各式中,  $W_f, W_i, W_c, W_o$  表示权重矩阵,  $b_f, b_i, b_c, b_o$  为偏置项,  $\sigma$  是权重函数,  $\tanh$  为激活函数。

### 2.3. EMD-LSTM 组合预测

在本节中将利用经验模型对原始数据进行经验分解, 并通过分解后所得的 IMF 和残差形成了基于 EMD-LSTM 神经网络预测模型。本文首先通过 EMD 算法将进站客流时间序列分解为  $n$  个分量, 然后对分量客流分别进行预测, 最终将各分量预测结果进行叠加得到最终预测值, 构建的短期客流预测方案如图 2 所示。

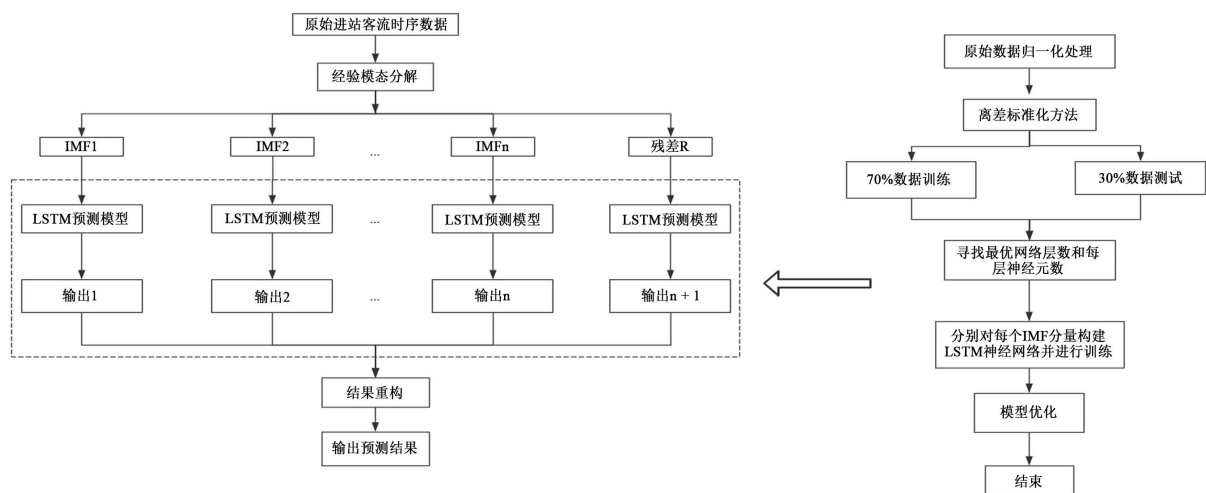


Figure 2. EMD-LSTM neural network combination  
图 2. EMD-LSTM 神经网络组合

## 3. 实验与结果

### 3.1. 数据来源及处理

本文通过现场人工录入采集, 得到济南地铁 1 号线长清区某地铁站点连续七天的进站实际客流数据, 以 5 分钟为统计间隔构造全日进站客流时序数据, 例如, 时间 06:00:00 表示聚合包含在 06:00:00 和 06:04:59 之间收集的数据信息, 以此类推。原始数据共有 1594 个客流观测值的连续时间序列, 并利用 Python3.10.4 软件的绘图功能得到了该站点连续七天的进站流量折线图, 如图 3 所示。

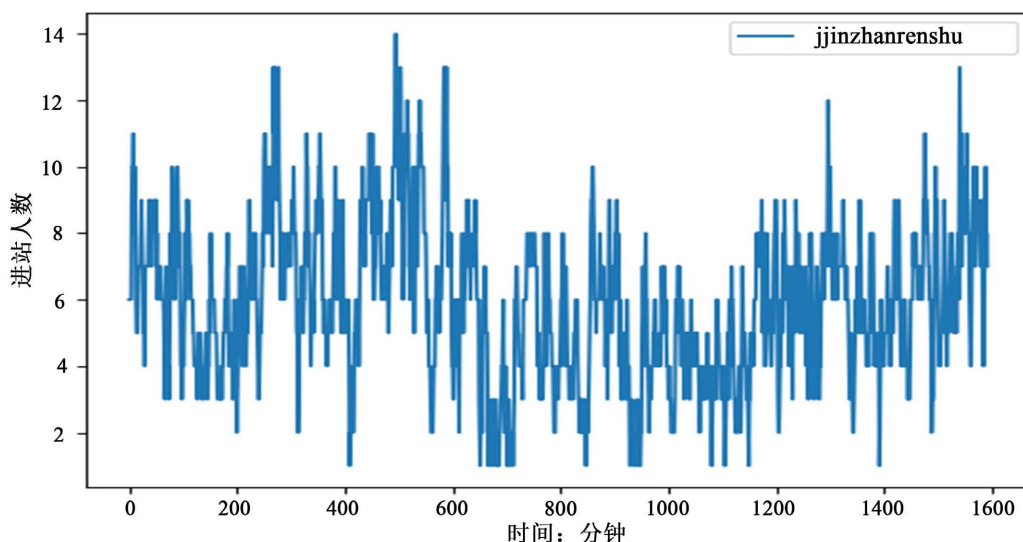


Figure 3. Line chart of inbound traffic

图 3. 进站流量折线图

从图 3 中可以发现, 在各个时段中都有几个波峰和波谷是不同的, 在序列中也有类似断崖式的变动。归一化处理是数据分析过程中必要的预处理操作, 这是为了在模型的训练过程中能更快地让参数趋于收敛, 同时对提高预测精度也有着重要的影响。因此在数据处理过程中, 经常会使用到离差标准化方法 (min-max 标准化方法)。离差标准化函数如式:

$$x^* = \frac{x - \min}{\max - \min} \quad (4)$$

其中,  $x$  为原始数据的任意一值, 而  $\min$  为数据集的最小值,  $\max$  则为数据集的最大值。  $x^*$  为经过离差标准化方法处理过后的新数据, 其值分布在  $[0, 1]$  之间。

### 3.2. 进站客流变化序列的 EMD 分解

经过对长宁区某站点的进站客流时间序列功率谱进行分析后可以发现, 日信号、早高峰信号、晚高峰信号构成了交通流量。利用 EMD 对交通流量序列进行分解, 共产生了如图 4 所示的 6 个 IMF ( $imf1 \sim imf6$ ), 以及 1 个 res。图中的 IMF1、IMF2、IMF3 的频率较高, 为进站数据中的高频分量, IMF4 开始序列的频率降低, 数据的波动性也开始降低, 为进站数据中的低频分量。使用经验模态(EMD)对连续 7 天的数据进行分解, 结果显示 EMD 方法能够较完整地解释了进站数据变动的常规周期项, 得出的 IMF 分量比原始数据更加稳定, 所提供的长期变化趋势也更加合理。

### 3.3. 网络评价指标

为了判定模型在预测城市轨道交通客流上的效果, 模拟评价指标可以判定地铁短时客流预测模型的误差。均方根偏差(RMSE)、平均绝对误差(MAE)、平均值绝对百分误差(MAPE)三种误差评价指标, 都可以对进站客流数据的预测准确度进行精准的评价。三种误差评价方法的具体式所示:

1) 均方根偏差(RMSE)。其作为衡量精确度的标准, 公式如下:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2} \quad (5)$$

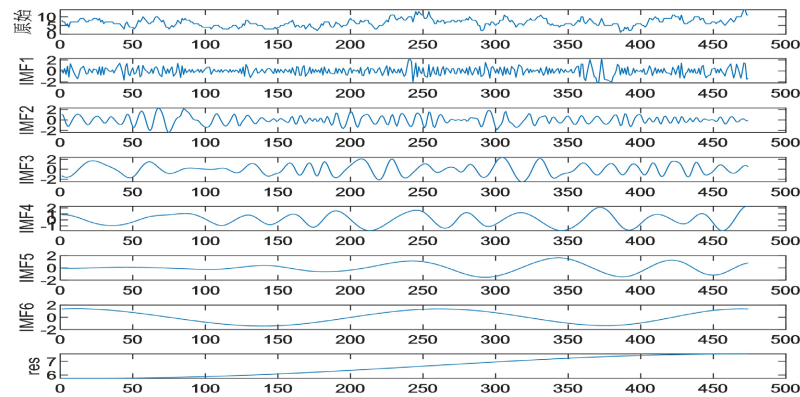


Figure 4. EMD decomposition results

图 4. EMD 分解结果

2) 平均绝对误差(MAE)。能更好地反映出预测值的误差情况,公式如下:

$$MAE = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{x}_i| \quad (6)$$

3) 平均值绝对百分误差(MAPE)。其值范围为 $[0, +\infty]$ , MAPE 值越小,代表模型的预测结果和实际值之间的偏差越小,则模型效果越好。公式如下:

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{x_i - \hat{x}_i}{x_i} \right| \times 100\% \quad (7)$$

式中,  $x_i$  为预测值,  $\hat{x}_i$  为真实值,  $i$  为数据集中该数据所处的位置数,为原始数据集中的真实值,是所有真实值的均值,  $n$  为原始数据集中的数据个数。

### 3.4. 实验结果

经 EMD 经验模态分解后所得的各 IMF 分量和残差,将其视为 LSTM 长短期记忆神经网络的输入向量。为了使实验结果更加准确,将上述分量进行区别后,取前 70% 作为训练集,后 30% 作为测试集,然后进行 LSTM 神经网络训练与预测,并将预测结果反归一化输出。对组合模型预测后得到的分量与预测结果合并,从而得出最终的预测结果。预测结果如图 5、图 6 所示。

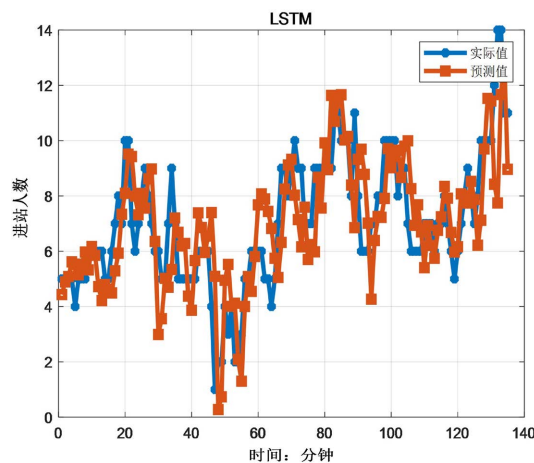


Figure 5. LSTM model results

图 5. LSTM 模型结果

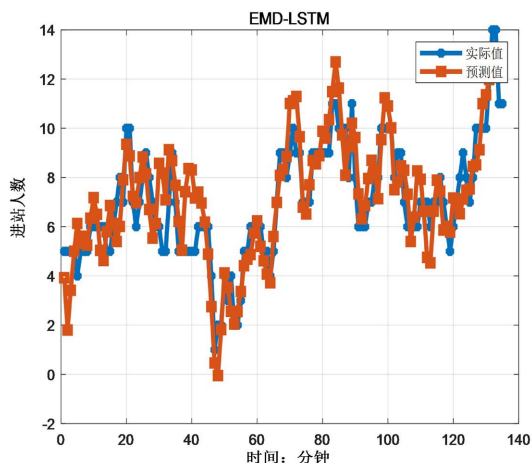


Figure 6. EMD-LSTM model results

图 6. EMD-LSTM 模型结果

将使用 EMD-LSTM 组合的神经网络方法和只使用 LSTM 模型的预测结果进行对比, 结果如图 7 所示。可以看出, EMD-LSTM 神经网络的预测结果线更接近于实际值, EMD-LSTM 神经预测能更好地模拟出数据的预测结果。说明本文所构建的模型具有更好的预测效果以及准确性。

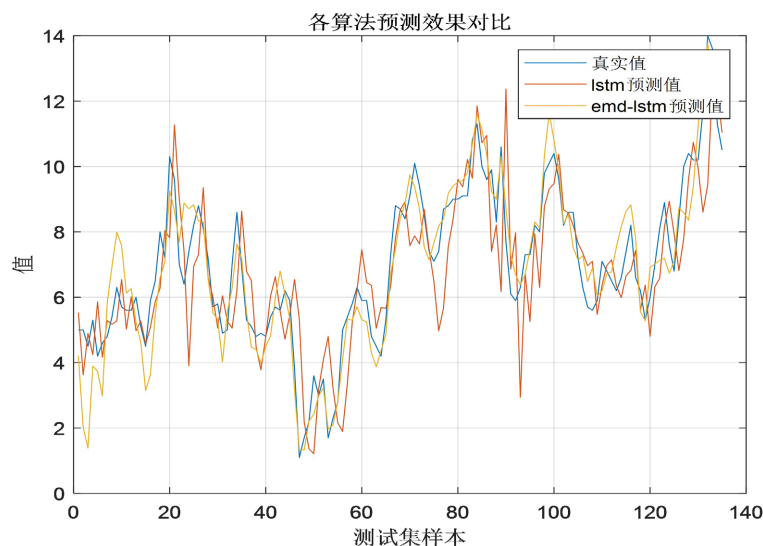


Figure 7. Comparison of prediction results between two models

图 7. 两种模型预测结果对比图

本文分别使用了均方根偏差(RMSE)、平均绝对误差(MAE)、平均值绝对百分误差(MAPE)三种误差评价指标对单独的 LSTM 预测模型和 EMD-LSTM 组合预测模型这两种模型进行了误差分析。其结果如表 1 所示。实验表明, LSTM 神经网络在处理不平稳性时间序列数据时具备强大的拟合和泛化能力, 同时发现了改进的空间, 即通过对原始数据进行 EMD 分解, 成功降低了波动性对预测精确度的影响, 从而显著提高了预测的准确性。

#### 4. 结语

随着城市轨道交通的不断发展, 公共交通的使用率不断上升, 拥堵问题时常出现。为了有效预测

**Table 1.** Evaluation and comparison of forecast results**表 1.** 预测结果评价对比

评价方式	LSTM	EMD-LSTM
均方根误差(RMSE)	1.6843	1.2337
平均绝对误差(MAE)	1.2914	0.99307
平均绝对百分比误差(MAPE)	17.2245%	12.6733%

进站客流, 进一步提升城市轨道交通的服务水平, 本文使用 EMD 方法与 LSTM 模型相结合, 由于进站数据流量分布的高度不稳定性, 引入了经验模态分解(EMD)算法, 将原始进站数据分解为多个固定分量, 实现了降噪和提高数据的鲁棒性, 以减少模型的干扰。接下来, 利用 LSTM 对每个独立的分量进行学习和预测, 最后将它们的预测结果叠加以获得最终的预测值。同时为了比较 LSTM 模型和 EMD-LSTM 模型的预测性能, 进行了模型效果检验。实际数据预测结果显示, 与 LSTM 相比, EMD-LSTM 组合模型的平均绝对误差(MAE)降低了 0.432, 均方根误差(RMSE)降低了 0.571, 平均百分比误差(MAPE)降低了 5.85 个百分点。综上所述, EMD-LSTM 组合模型的预测精度明显优于单独的 LSTM 模型。这种方法在进站客流的短期预测方面表现出色, 并且在其他时间序列预测领域也有广泛的应用前景。

## 参考文献

- [1] Huang, N.E., Shen, Z., Long, S.R., *et al.* (1998) The Empirical Mode Decomposition and the Hilbert Spectrum for Nonlinear and Non-Stationary Time Series Analysis. *Proceedings: Mathematical, Physical and Engineering Sciences*, **454**, 903-995. <https://doi.org/10.1098/rspa.1998.0193>
- [2] Wu, Z. and Huang, N.E. (2009) Ensemble Empirical Mode Decomposition: A Noise-Assisted Data Analysis Method. *Advances in Adaptive Data Analysis*, **1**, 1-41. <https://doi.org/10.1142/S1793536909000047>
- [3] 曾阳艳, 苏雅, 等. 基于 EMD-LSTM 神经网络的交通流量预测模型[J]. 湖南商学院学报, 2021, 28(3): 110-115.
- [4] Huang, N.E., Shih, H.H., Shen, Z., *et al.* (2000) The Ages of Large Amplitude Coastal Seiches on the Caribbean Coast of Puerto Rico. *Journal of Physical Oceanography*, **30**, 2001-2012. [https://doi.org/10.1175/1520-0485\(2000\)030<2001:TAOLAC>2.0.CO;2](https://doi.org/10.1175/1520-0485(2000)030<2001:TAOLAC>2.0.CO;2)
- [5] Huang, N.E., Chern, C.C., Huang, K., *et al.* (2001) A New Spectral Representation of Earthquake Data: Hilbert Spectral Analysis of Station TCU129, Chi-Chi, Taiwan, 21 September 1999. *Bulletin of the Seismological Society of America*, **91**, 1310-1338. <https://doi.org/10.1785/0120000735>
- [6] 程军圣, 于德介, 等. EMD 方法在转子局部碰摩故障诊断中的应用[J]. 振动、测试与诊断, 2006, 26(1): 24-27.
- [7] 陈志全, 雷景生. 基于 MPSO-LSTM 模型的股票指数预测研究[J]. 现代信息科技, 2021, 5(4): 1-4.
- [8] 王振龙, 顾岚. 时间序列分析[M]. 北京: 中国统计出版社, 2000.
- [9] 罗凤曼. 时间序列预测模型及其算法研究[D]: [硕士学位论文]. 成都: 四川大学, 2020.