

# 山东财经大学

# 本科毕业论文(设计)

题目：基于时序分析模型的金融数据变化趋势预测

学    院 计算机科学与技术学院  
专    业 计算机科学与技术专业  
班    级 计科 1903 班  
学    号 20191334109  
姓    名 李成龙  
指导教师 刘    慧

山东财经大学教务处制

二〇二三年五月

## 山东财经大学学士学位论文原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不含任何其他个人或集体已经发表或撰写过的研究成果。对本文的研究做出重要贡献的个人和集体，均已在论文中作了明确的说明并表示了谢意。本声明的法律结果由本人承担。

学位论文作者签名：\_\_\_\_\_

\_\_\_\_\_年\_\_\_\_月\_\_\_\_日

---

## 山东财经大学关于论文使用授权的说明

本人完全了解山东财经大学有关保留、使用学士学位论文的规定，即：学校有权保留、送交论文的复印件，允许论文被查阅，学校可以公布论文的全部或部分内容，可以采用影印或其他复制手段保存论文。

指导教师签名：\_\_\_\_\_

\_\_\_\_\_年\_\_\_\_月\_\_\_\_日

论文作者签名：\_\_\_\_\_

\_\_\_\_\_年\_\_\_\_月\_\_\_\_日

# 基于时序分析模型的金融数据变化趋势预测

## 摘 要

金融市场反映着经济状态，但由于其高度复杂性，预测其变化趋势具有挑战性。神经网络和机器学习方法可应用于金融数据预测，帮助投资者把握市场动态，规避风险，制定更为合理的投资策略，为金融监管提供更有效的政策。技术进步提升了预测算法和模型的可靠性，构建精度高的预测模型可帮助分析未来价格走势，制定更符合金融经济发展方向的政策。

本文研究了基于差分整合移动平均自回归模型(Autoregressive Integrated Moving Average model, ARIMA)的股票价格预测方法、基于 LSTM 神经网络的股票价格预测方法，包括基于长短期记忆网络(Long Short Time Memory, LSTM)和结合注意机制的长短期记忆网络(Attention Long Short Time Memory, AT-LSTM)模型的股票价格预测方法。针对 ARIMA、LSTM 和 AT-LSTM 的基本原理进行介绍，并将其在中国平安股票、贵州茅台股票、中国石油股票和中国工商股票四个股票数据集上进行了验证。

本文的股票价格预测方法基于 Matlab 和 PyTorch 框架，通过计算三种模型的 MAE、MSE 和 RMSE 评价指标对模型的预测效果进行评估并分析其预测的准确度。实验结论如下：AT-LSTM 模型在中国平安股票、中国石油股票和中国工商股票三种股票数据集上的预测效果最好，由于贵州茅台股票数据的波动性较大、周期性不明显，导致 ARIMA、LSTM、AT-LSTM 三种模型在该数据集上的预测效果都较差。因此，在金融数据预测任务上，要根据股票数据集具体的数据趋势、数据特征及预测时间序列的长短综合考虑选择合适的预测模型，以达到更好的预测效果。本文研究的基于 LSTM 模型的股票价格预测方法及实验结果对股票价格预测及其金融风险识别的解决及研究提供了重要的参考依据。

**关键词：**股票预测；ARIMA；LSTM；注意力机制

## ABSTRACT

The financial market reflects the economic state, but predicting its trends is challenging due to its high complexity. Neural networks and machine learning methods can be applied to financial data forecasting to help investors grasp market dynamics, mitigate risks, formulate more reasonable investment strategies, and provide more effective policies for financial regulation. Technological advancements have improved the reliability of prediction algorithms and models, and building accurate prediction models can help analyze future price trends and formulate policies that are more in line with the direction of financial and economic development.

This study investigates stock price prediction methods based on the Autoregressive Integrated Moving Average model (ARIMA), the Long Short Time Memory (LSTM) neural network, and the Attention Long Short Time Memory (AT-LSTM) model. The basic principles of ARIMA, LSTM, and AT-LSTM are introduced, and their performance is verified on four stock datasets of China Ping An, Guizhou Maotai, China Petroleum, and China Industrial and Commercial Bank.

The stock price prediction methods in this study are based on the Matlab and PyTorch frameworks, and the MAE, MSE, and RMSE evaluation indicators are used to evaluate and analyze the prediction accuracy of the three models. The experimental conclusions are as follows: the AT-LSTM model has the best prediction performance on the three stock datasets of China Ping An, China Petroleum, and China Industrial and Commercial Bank. Due to the high volatility and unclear periodicity of Guizhou Maotai stock data, the prediction performance of the three models is poor on this dataset. Therefore, when it comes to financial data prediction tasks, it is necessary to select a suitable prediction model based on the specific data trends, data characteristics, and the length of the predicted time series to achieve better prediction performance. The stock price prediction method based on the LSTM model and the experimental results of this study provide important reference for stock price prediction and financial risk identification research and solutions.

**Keywords:** Stock Prediction; ARIMA; LSTM; Attention Mechanism

## 目录

一、绪论.....	1
(一) 课题研究背景与意义.....	1
(二) 国内外研究现状.....	1
1.国外研究现状.....	1
2.国内研究现状.....	2
3.国内外研究文献综述简介.....	2
(三) 主要研究内容.....	3
(四) 论文组织结构.....	3
二、金融时序分析相关理论.....	4
(一) 金融时间序列简介.....	4
1.金融时间序列介绍.....	4
2.金融时间序列预测原理.....	4
(二) 金融时间序列数据预处理.....	5
1.输入特征选取.....	5
2.缺失数据处理.....	5
3.相关系数分析.....	6
4.相关评价指标.....	6
(三) ARIMA 模型的基本理论.....	7
1.ARIMA 原理.....	7
2.ARIMA 模型公式推导.....	7
(四) LSTM 模型的基本理论.....	8
1.LSTM 原理.....	9
2.LSTM 模型公式推导.....	9
(五) AT-LSTM 模型的基本理论.....	10
1.AT-LSTM 原理.....	10
2.AT-LSTM 模型公式推导.....	11
三、数据集介绍.....	12
四、基于 ARIMA 模型的金融数据变化趋势预测.....	16
(一) 数据预处理.....	16
1.数据平稳性检验.....	16
2.数据平稳化.....	17
(二) 实验环境及参数确定.....	18
1.实验环境.....	18
2.参数确定.....	19
(三) 实验结果分析.....	20
五、基于 LSTM 和 AT-LSTM 模型的金融数据变化趋势预测.....	22
(一) 数据预处理.....	22
(二) 基于 LSTM 模型的实证分析.....	23
1.实验环境.....	23
2.模型搭建.....	23
3.实验结果分析.....	24
(三) 基于 AT-LSTM 模型的实证分析.....	25

1.模型搭建.....	25
2.实验结果分析.....	26
六、总结与展望.....	28
（一）研究总结.....	28
（二）研究展望.....	28
参考文献.....	30
致谢.....	31

## 一、绪论

### （一）课题研究背景与意义

金融作为经济社会运转中不可或缺的组成部分，其市场表现直接反映了一个国家的发展状态。随着全球化进程的不断加速，任何一个地区的经济变动都可能引起全球金融市场的波动，因此金融市场的稳定性变得尤为重要。然而，由于金融市场的高度复杂性，金融时间序列数据一直以来都具有非线性、高噪声和非平稳性等特征，因此预测其变化趋势具有极大的挑战性。因此，研究如何利用神经网络和机器学习方法对金融数据进行预测，以实现对其变化趋势的精准预测，具有极其重要的学术意义和实践价值。

同时，随着计算机技术和数据挖掘方法的快速发展，神经网络和机器学习技术已广泛应用于金融市场预测领域。通过运用这些技术，我们可以有效地处理金融市场的非线性、高噪声和非平稳性等特征，以精确预测其变化趋势。例如，利用长短期记忆网络、自回归整合移动平均模型等方法可以实现对金融数据变化趋势的预测。

金融数据变化趋势的预测对于投资者、机构投资者和金融监管机构等各方具有重要意义。首先，金融市场的不确定性使得其变化与多种因素有关，而准确地预测金融市场的变化趋势有助于这些利益相关方更好地应对不确定性。其次，金融数据变化趋势的预测有助于投资者把握市场动态，制定更为合理的投资策略，从而获得更高、更稳定的回报。对于机构投资者和金融监管机构来说，准确的预测也有助于规避投资风险和制定更有效的政策。最后，随着数据挖掘、机器学习和人工智能等技术的不断进步，金融数据变化趋势的算法和模型得到了极大的提升，预测结果更加准确、稳定和可靠，为金融市场的决策提供了更好的技术支持。

因此，我们可以利用现有大量高维度的金融数据，构建高精度的金融预测模型，以从中获得有用的信息。这将有助于投资者分析金融数据未来的价格走势，以及帮助金融监管者制定更符合金融经济发展方向的政策。随着技术的不断提升和数据的不断积累，金融预测模型的精度和可靠性将不断得到提高，为金融市场的稳定运转提供更好的保障。

### （二）国内外研究现状

#### 1. 国外研究现状

1997年 Sepp Hochreiter 和 Jurgen Schmidhuber<sup>[1]</sup>提出了 LSTM 神经网络模型，通过引入一种新的、高效的、基于梯度的方法，提高了学习速度，能够解决复杂的、时间序列方面的问题，该模型现在广泛应用于自然语言处理、金融时序预测等领域。2002年，Francis E.H. Tay 和 L.J. Cao<sup>[2]</sup>提出了一种基于 SVM 的改进模型，即 C-级联递增支持向量机，通过修改支持向量机中的正则化风险函数，使得近期不敏感的错误比远期不敏感的错误受到更多的惩罚，从而实现对非平稳的金融时间序列模拟。此外，在 2012 年，Md. Rafiul Hassan 等人<sup>[3]</sup>提出了一种基于 HMM、模糊逻辑和多目标进化（EA）的混合算法，该算法通过使用 HMM 对每个数据模式的对数似然分数来对数据进行排序，并使用排序的数据生成模糊

规则，从而建立模糊模型来对非线性时间序列数据进行预测。

2015 年，Javad Zahedi 和 Mohammad Mahdi Rounaghi<sup>[4]</sup>将人工神经网络（ANNs）和主成分分析（PCA）方法结合起来，提出了一种用于评估德黑兰证券交易所价格可预测性的模型。该模型使用了 20 个会计变量来评估证券交易所价格的可预测性，并通过利用主成分分析的拟合度来确定价格的有效因素。在 2017 年，Salim Lahmiri<sup>[5]</sup>提出了一种基于奇异谱分析（SSA）、支持向量回归（SVR）和粒子群优化（PSO）的模型，通过与小波变换（WT）和前馈神经网络（FFN）融合的模型进行对比，发现 SSA-PSO-SVR 模型比较适用于高噪声金融时间序列的分析和预测。2019 年，Taewook Kim 和 Ha Young Kim<sup>[6]</sup>提出了一种 LSTM-CNN 的融合模型，该模型结合了从同一数据的不同表示中学习到的特征，即股票时间序列和股票图标图像，以预测股票价格。他们发现使用同一数据的时间和图像特征的组合可以有效地减少预测误差。2020 年，Shihua Luo 和 Cong Tian<sup>[7]</sup>提出了一种快速分步网格搜索方法（SGS），用于优化 LSTM 网络的超参数，从而获得更高的效率和更小的均方根误差。

## 2. 国内研究现状

2008 年，裴双喜<sup>[8]</sup>使用 ARMA 和 ARCH 模型利用上证指数，对 ARMA-ARCH 模型进行建模分析，证明了其具有良好的短期预测效果。2010 年，张文霄<sup>[9]</sup>将 PSO 算法应用于 BP 神经网络的优化中，建立了基于 BP 神经网络的股价预测模型。2015 年，孙瑞奇<sup>[10]</sup>利用拟牛顿法原理对 LSTM 神经网络中的学习率进行改进，实现了 LSTM 的改进。2017 年，王谨平<sup>[11]</sup>在误差计算中将历史数据按照与当前时间的远近赋予相应的权值，同时加入时序数据随机过程，改进了 Elman 神经网络时间序列预测模型（GT-Elman），增强了其对时间序列的预测性能。2018 年，周凌寒<sup>[12]</sup>利用 LSTM 神经网络并融合基本情感特征，以提高模型预测准确性。2019 年，陈璐<sup>[13]</sup>提出了一种由自适应噪声的完整集成经验模态分解并添加了注意力机制的长短期记忆网络，并采用上证 50 指数对模型进行验证，证明该模型比其他模型具有更好的预测误差。2020 年，赵薇<sup>[14]</sup>提出了一种融合了 ARIMA 模型和 PSO 优化算法的 ARIMA-PSO-LSTM 模型，并对 PSO-LSTM 和 LSTM 的实验结果进行分析，得出提出的模型预测效果更好的结论。2022 年，李庆涛<sup>[16]</sup>提出了基于板块效应的 CNN-LSTM 深度学习混合模型，解决了单一神经网络模型预测精度差的问题，并提出了一种基于多元特征的生成对抗网络模型。

## 3. 国内外研究文献综述简介

国内外对于金融时间序列研究较多，国外学者主要注重于对模型的创造和改进，国内学者主要对模型进行融合、改进和应用，并通过国内金融市场对改进模型的实证分析，得出改进模型对国内金融市场经济数据预测效果较好的结论。

随着计算机的算力逐渐提高，大量的算法应用到金融时间序列数据的预测方向上，比如传统的统计学模型，例如 ARIMA 模型、CARCH 模型等，到机器学习和深度学习，比如 SVM、LSTM、MLP、HMM 等，单一的神经网络模型对金融数据的预测取得较大的成功，



但是也有其缺陷，比如参数较少导致预测的精度较差，故国内外大量的学者对模型进行融合，不断在基础的单一神经网络模型上进行改进，比如将 PSO 算法放入模型的参数优化中，利用拟牛顿法原理对 LSTM 中学习率的改进，再有利用模糊学习、强化学、小波理论等和 LSTM 模型相结合以提高模型预测金融时间序列数据的准确性。目前单一的模型已不再是研究的方向，大模型、组合模型已经成为金融数据预测问题上的热门方向，也是国内外学者和从业人员越来越关注的技术发展方向。

### （三）主要研究内容

本文旨在探讨时间序列模型在金融数据预测方面的研究历程、问题以及现阶段的解决方案。首先介绍了时间序列模型的研究历史和应用场景，接着分析了目前时序数据预测存在的问题，主要包括数据噪声、非线性、非稳定性等方面。然后本文在 ARIMA、LSTM 和 AT-LSTM 三种时间序列模型基础上进行股票数据集收盘价数据预测的任务，并在中国平安股票(000001.SZ)、贵州茅台股票(600519.SS)、中国石油股票(601857.SS)和中国工商银行股票(601398.SS)四个股票数据集上进行了验证。使用 MAE、MSE 和 RMSE 评价指标对三个模型的预测效果进行评估并分析其预测的准确度。

最后，本文对这三种模型进行了优缺点的总结，并指出了未来时间序列模型研究的方向。其中，ARIMA 模型具有简单易用、可解释性强的优点，但是对于非线性和非稳定的数据表现不佳；LSTM 模型可以处理非线性和长期依赖的数据，但是对于噪声和非平稳的数据存在一定的问题；AT-LSTM 模型相比于 LSTM 模型更具有鲁棒性和可解释性，但是计算量较大，需要更长的训练时间。实验结果显示，AT-LSTM 模型在中国平安股票、中国石油股票和中国工商银行股票三种股票数据集上的预测效果最好，但由于贵州茅台股票数据的波动性较大、周期性不明显，导致 ARIMA、LSTM、AT-LSTM 三种模型在该数据集上的预测效果都较差。因此，在金融数据预测任务上，要根据股票数据集具体的数据趋势、数据特征及预测时间序列的长短综合考虑选择合适的预测模型，以达到更好的预测效果。

### （四）论文组织结构

第一章：绪论。详细提出了金融时间序列预测的重要作用和目前遇到的问题，通过介绍国内外学者对于金融时序分析预测的研究成果，确定未来的研究方向，并结合国内外学者发表的论文，总结出近年来的金融时间序列数据预测的主要研究成果。最后写出了本文的主要研究内容和文章结构安排。

第二章：金融时序分析相关理论。该章节首先介绍了金融时序的定义，特征及对金融时间序列预测的原理。通过对金融时间序列的介绍引出对金融时间序列数据的处理，包括了输入特征的选取，缺失数据的处理，皮尔逊相关系数的分析，最后得出预测数据后对模型进行相关评价指标的分析。接下来阐述 ARIMA、LSTM 和 AT-LSTM 模型的原理，并为下面单一神经网络模型和复杂神经网络模型的实证分析提供理论基础。

第三章：数据集介绍。该章节介绍了相关数据集，并分析其统计学特征。通过对其数

据的分析增加实验效率。

第四章：基于 ARIMA 模型的金融数据变化趋势的预测。本章围绕着 ARIMA 进行阐述金融时间序列数据预测，详细论述了数据的来源和数据的分析与预处理操作，选取了平安银行股价的相关数据并利用 ARIMA 模型的基本原理和训练的过程，给出 ARIMA 模型的训练结果和预测结果的分析。

第五章：基于 LSTM 神经网络模型的金融数据变化趋势的预测。本章围绕着 LSTM 及 AT-LSTM 神经网络模型对平安银行股价进行实证分析，详细介绍其原理和训练过程，并给出预测结果的评价指标并对模型的预测结果进行评价和分析。

第六章：总结与展望。总结本文阐述的内容，并通过分析文中提到的模型，来预测金融数据的变化趋势，得出结论，并对未来金融时间序列的相关研究给出方向。对本文提到的四个股票价格预测模型进行总结对比分析，分析模型的优缺点，提出未来需要改进的方向。

## 二、金融时序分析相关理论

### （一）金融时间序列简介

#### 1. 金融时间序列介绍

时间序列是一组按照时间顺序排列的数据点，这些数据点通常按照一定时间间隔采集的观测值或测量值。时间序列的数据通常用于研究时间相关的趋势、季节性变化、周期性变化等。时间序列可以表示为一个由一系列时间点和相应的观测值组成的序列，其数学表示为  $Y_1, Y_2, \dots, Y_t$  其中， $Y_t$  代表在时间点  $t$  处观测到的值。时间点  $t$  通常是等间隔的、例如每天、每周或每月。

金融时间序列是指按照时间顺序记录的各种金融变量数据，如股票价格、汇率、利率、成交量等。金融时间序列具有以下的特点：

非常规分布：金融时间序列通常不符合正态分布，具有尖峰厚尾的特征，即存在大量极端值和异常点。

长期依赖：金融时间序列具有长期的相关性，即过去的信息对未来的预测具有重要的影响。

非稳定性：金融时间序列通常不具有稳定性，即统计特性会随时间而变化。

季节性和周期性：金融时间序列中存在明显的季节性和周期性变化。

聚集效应：金融时间序列中具有聚集效应，即大的波动通常会伴随着更多的波动。

#### 2. 金融时间序列预测原理

分析金融市场的变化趋势从本质上来说就是根据历史的金融时间序列数据的走势预测出未来的金融市场价格的变化趋势。Fama<sup>[17]</sup>提出了有效市场其指的是在合格市场中有许多理性且以利润最大化为目标的、积极竞争的交易者，他想要通过当前的可用的信息来预

测未来市场中资产的价值。因此通过市场有效性假说可以得出金融市场的数据是可以预测的。

金融时间序列的预测一般分为三种<sup>[18]</sup>，第一是基于计量经济学模型的金融时间序列数据预测，用到了许多模型，比如 AR、MA、ARMA 等模型，其模型都适用于平稳的时间序列数据，但在处理非平稳时间序列的数据方面存在一定的局限性，因此有人提出了 ARIMA 模型，来处理非平稳时间序列的数据，并对数据进行预测。

第二是基于机器学习算法的金融时间序列数据的预测，可以利用 BP 神经网络、SVM（支持向量机）、随机森林等机器学习的算法，实现了一种非参数化、非线性、数据驱动模型，机器学习模型可以很好的处理金融时间序列数据的非平稳性、非线性等特性，更容易挖掘金融时间序列数据之间的关系。因此机器学习模型相对于传统的计量经济学模型对于金融时间序列数据预测更具有优势。

最后是基于深度学习算法的金融时间序列数据预测，比如 CNN、RNN、LSTM 等神经网络或者组合神经网络模型，相对于机器学习，深度学习对金融时间序列数据的特征提取更加有效，可以更好的刻画金融时间序列数据的相关特征。所以相对于前两种预测模型，基于深度学习模型的预测更具有优势。

## （二）金融时间序列数据预处理

数据的预处理是非常重要的，因为金融时间序列数据具有高噪声、不平稳性和非常规性等特点，导致其数据是非常复杂的，因此正确的进行数据预处理是必不可少的。根据金融时间序列数据的特征，采用恰当的数据预处理可以有效地降低时间序列预测的误差，同时有效的特征选取和获取的算法可以应用于金融时间序列，通过特征选取的算法选择最佳的输入特征向量、减少数据的维度和缩短训练时的时间可以减少算力消耗。数据预处理可以较好的选择输入的变量，划分出训练集和测试集，筛选出较好的向量，并检测出异常的数据根据相应的缺失数据处理机制来进行数据的填充，同时对数据输入特征向量进行分解，降低预测的误差。因此数据预处理机制是必不可少的。

### 1. 输入特征选取

金融时间序列数据中，因其数据具有很高的维度，且噪声较多，如果不对其进行输入特征的选取<sup>[19]</sup>，预测效果会交叉，因此我们对数据的选择需要尽可能地考虑到影响股票价格的特征因素。比较基础的有开盘价、最高价、最低价、收盘价等。其中根据技术面分析股票价格的时候还有技术指标，可以通过股票的基本面来判断股票价格的基本走势。比如股票均线、异同移动平均线、布林带等。

### 2. 缺失数据处理

金融市场的时序数据中，股票存在休市或者停牌等状况，所以金融时间序列数据，即股票的数据并非是完全连续的时间序列数据，因此获取的数据往往是缺失的或者异常的，

因此处理缺失数据是必要的，下面介绍几种处理缺失数据的方法：

**删除缺失数据：**直接删除缺失数据的日期，在金融时间序列预测过程中选取得数据往往是以年记，删除几天的数据并不会对整体的趋势有影响，因此可以选择直接删除缺失数据。

**填充缺失数据：**可以使用统计方式的填充，比如就近填充、特征值填充、线性插值等填充方法，实现缺失数据的再现。

### 3. 相关系数分析

股票的价格收到了很多因素的影响，但是不同的因素对股票价格的影响程度不一样，因此可以利用相关系数来选取相关性较好的特征来对股票价格进行预测，可以提高股票价格预测的精度。不同股票的特征之间对股票价格的影响程度不同，因此可以利用相关性差异，选取与股票价格强相关的特征进行数据预测。本文利用皮尔逊相关系数来衡量因素与股票价格的相关性<sup>[20]</sup>。

皮尔逊相关系数是用于度量两个变量  $X$  和  $Y$  之间的相关（线性相关），其值介于-1 到 1 之间。其计算公式如下：

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

两个向量  $X$  和  $Y$ ，计算出的皮尔逊相关系数含义如下：

- (1) 当相关系数为 0 时， $X$  和  $Y$  两向量不相关。
- (2) 当  $X$  的值增大（减小）， $Y$  值减小（增大）， $X$  和  $Y$  两个向量负相关，相关系数在 -1.0 到 0.0 之间。
- (3) 当  $X$  的值增大（减小）， $Y$  值增大（减小）， $X$  和  $Y$  两个向量正相关，相关系数在 0.0 到 +1.0 之间。

### 4. 相关评价指标

为了评价模型对股票收盘价预测的效果，本文选取平均绝对误差（Mean Absolute Error, MAE）、均方误差（Mean Squared Error, MSE）、均方根误差（Root Mean Squared Error, RMSE）三种评价指标。

MAE 指模型预测值  $f(x)$  与样本真实值  $y$  之间距离的平均值。其公式为：

$$MAE = \frac{1}{m} \sum_{i=1}^m |y_i - f(x_i)| \quad (2-2)$$

MSE 指预测值与真实值的匹配程度。其公式为：

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (2-3)$$

RMSE 是均方根误差，也称标准误差，即均分误差的算数平方根。均分误差的量纲与数据量纲不同，不能直接反映离散程度，故在均方误差上开平方根。其公式为：

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (2-4)$$

### （三）ARIMA 模型的基本理论

差分自回归移动平均（ARIMA）模型是一种用于时间序列分析的统计模型<sup>[21]</sup>，它可以预测时间序列数据的未来走势。ARIMA 模型基于时间序列数据的自相关性、趋势性和季节性，通过差分、自回归和移动平均的组合来描述时间序列的性质。

#### 1. ARIMA 原理

ARIMA 模型可以分为三个部分：AR、I 和 MA。其中，AR 代表自回归模型，I 代表差分模型，MA 代表移动平均模型。ARIMA(p,d,q)模型的意义如下：

（1）p 表示自回归模型（AR）中使用的滞后阶数，即时间序列数据自身的历史数据对当前值的影响；

（2）d 表示差分阶数，即对时间序列进行差分的次数，用于平稳化时间序列数据；

（3）q 表示移动平均模型（MA）中使用的滞后阶数，即时间序列数据误差的历史数据对当前值的影响。

#### 2. ARIMA 模型公式推导

ARIMA 模型的公式如下：

$$ARIMA(p,d,q): y'_t = c + \sum_{i=1}^p \phi_i y'_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} + \epsilon_t \quad (2-5)$$

其中， $y'_t$  是经过  $d$  阶差分后的时间序列数据， $c$  是常数， $\phi_i$  和  $\theta_j$  是 AR 和 MA 模型中的系数， $\epsilon_t$  是误差项。

AR 模型用于描述时间序列数据的自回归性质，即当前值受到过去  $p$  个时间点的历史值的影响，可以表示为：

$$AR(p): y'_t = c + \sum_{i=1}^p \phi_i y'_{t-i} + \epsilon_t \quad (2-6)$$

其中， $\phi_i$  是 AR 模型中的系数。

MA 模型用于描述时间序列数据的移动平均性质，即当前值受到过去  $q$  个时间点的误差的影响，可以表示为：

$$MA(q): y'_t = c + \sum_{i=1}^q \phi_i y'_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} + \epsilon_t \quad (2-7)$$

ARIMA 模型就是在 ARMA 模型的基础上，增加了差分的操作，用于平稳化时间序列数据。具体来说，对于非平稳时间序列数据，可以通过差分操作将其转化为平稳时间序列

数据，然后再应用 ARMA 模型进行建模和预测。故 ARIMA 模型的推导公式如下：

首先，对于一个时间序列  $y_t$ ，定义其  $d$  阶差分为：

$$y'_t = \nabla^d y_t = (1-L)^d y_t \quad (2-8)$$

其中， $L$  是向后移动一个时间步的滞后算子， $(1-L)$  表示对时间序列进行一阶差分操作， $(1-L)^2$  表示对时间序列进行二阶差分操作，以此类推。

接下来，将  $y'_t$  带入到 ARMA(p,q) 模型中，得到：

$$y'_t = c + \sum_{i=1}^p \phi_i y'_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} + \epsilon_t \quad (2-9)$$

将  $y'_t$  带入上式，得到 ARIMA(p,d,q) 模型的推导公式：

$$(1-L)^d y_t = c + \sum_{i=1}^p \phi_i (1-L)^d y_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} + \epsilon_t \quad (2-10)$$

展开上式，得到：

$$(1 - \sum_{k=1}^d C_d^k L^k) y_t = c + \sum_{i=1}^p \phi_i (1-L)^d y_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} + \epsilon_t \quad (2-11)$$

其中， $C_d^k$  是组合数。

对于 ARIMA(p,d,q) 模型，需要选择合适的参数(p,d,q)，使得模型能够准确的描述时间序列数据的性质，并能够进行准确的预测。可以使用统计方法，比如赤池信息准则和贝叶斯信息准则等指标，来进行模型参数的选择。

确定了 ARIMA 模型的参数(p,d,q)后，可以使用模型来预测时间序列的未来值。下面是 ARIMA 模型的预测公式：

$$\hat{y}_{t+h} = c + \sum_{i=1}^p \phi_i (1-L)^d y_{t+h-i} + \sum_{j=1}^q \theta_j \epsilon_{t+h-j} \quad (2-12)$$

其中， $\hat{y}_{t+h}$  是  $t+h$  时刻的预测值， $h$  表示预测的时间步长。在进行预测时，需要先确定预测的时间步长  $h$ ，然后使用已知的历史数据进行建模，得到 ARIMA(p,d,q) 模型的参数，最后带入上式进行预测。

需要注意的是，ARIMA 模型预测的精度受到多个因素的影响，如模型的参数选择、历史数据的质量、预测时间步长等。因此，在使用 ARIMA 模型进行预测时，需要进行充分的数据分析和模型评估，以确保预测结果的准确性和可靠性。

#### (四) LSTM 模型的基本理论

长短期记忆神经网络(LSTM)是一种常用于处理序列数据的循环神经网络架构，它通过增加一个记忆单元和三个门控单元来解决传统 RNN 存在的梯度消失和梯度爆炸问题，可用于金融时间序列数据的预测<sup>[22]</sup>。

## 1. LSTM 原理

LSTM 的主要思想是在每个时间步上，维护一个长期的状态向量  $C_t$ ，并使用三个门控单元来控制信息的流动。其中，输入门（Input Gate） $i_t$  决定哪些信息应该被添加到状态向量中，遗忘门（Forget Gate） $f_t$  决定哪些信息应该被丢弃，输出门（Output Gate） $o_t$  决定哪些信息应该被输出。LSTM 神经网络结构如图 2-1：

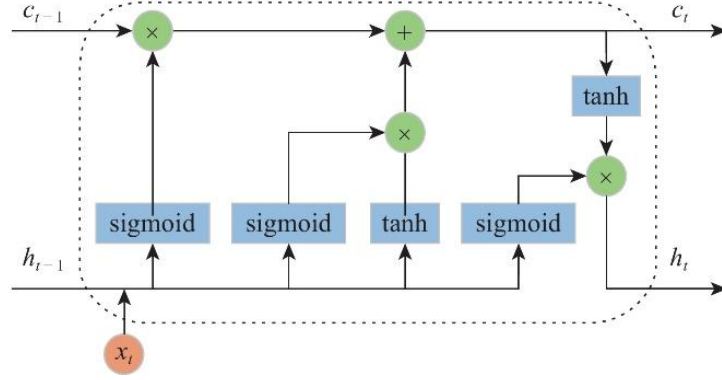


图 2-1 LSTM 神经网络结构图

在每个时间步  $t$ ，LSTM 的输入包括当前的输入向量  $x_t$ 、上一个时间步的输出向量  $h_{t-1}$  以及上一个时间步的状态向量  $C_{t-1}$ 。LSTM 的输出包括当前时间步的输出向量  $h_t$  和当前时间步的状态向量  $C_t$ 。

## 2. LSTM 模型公式推导

LSTM 的具体计算过程如下：

### （1）输入门 $i_t$

输入门  $i_t$  决定哪些信息应该被添加到状态向量  $C_t$  中。它首先根据当前的输入向量  $x_t$  和上一个时间步的输出向量  $h_{t-1}$  计算一个候选门控向量  $\tilde{C}_t$ ，然后使用 Sigmoid 函数将其转换为  $[0,1]$  范围内的值。即：

$$\tilde{C}_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (2-13)$$

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad (2-14)$$

其中， $W_{xc}$ 、 $W_{hc}$ 、 $b_c$  分别是输入向量、输出向量与候选门控向量的权重矩阵和偏置向量， $W_{xi}$ 、 $W_{hi}$ 、 $b_i$  分别是输入向量、输出向量与输入门控向量的权重矩阵和偏置向量， $\sigma$  是 Sigmoid 函数， $\tanh$  是双曲正切函数。

### （2）遗忘门 $f_t$

遗忘门  $f_t$  决定哪些信息应该被丢失。它首先根据当前的输入向量  $x_t$  和上一个时间步的输出向量  $h_{t-1}$  计算一个候选遗忘向量  $\tilde{f}_t$ ，然后使用 Sigmoid 函数将其转换在  $[0,1]$  范围内的

值。即：

$$\tilde{f}_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_f) \quad (2-15)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \quad (2-16)$$

其中， $W_{xf}$ 、 $W_{hf}$ 、 $b_f$  分别是输入向量、输出向量与遗忘门控向量的权重矩阵和偏置向量。

### (3) 状态更新 $C_t$

根据输入门  $i_t$  和遗忘门  $f_t$ ，可以得到当前时间步的状态更新向量  $C_t$ 。具体地，我们将  $\tilde{C}_t$  乘上输入门  $i_t$ ，再将上一个时间步地状态向量  $C_{t-1}$  乘上遗忘门  $f_t$ ，将它们相加即可得到  $C_t$ 。即：

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (2-17)$$

其中， $\odot$  表示逐元素相乘。

### (4) 输出门 $o_t$

输出门  $o_t$  决定哪些信息应该被输出。它首先根据当前的输入向量  $x_t$  和上一个时间步的输出向量  $h_{t-1}$  计算一个候选输出向量  $\tilde{h}_t$ ，然后使用 Sigmoid 函数将其转换为在  $[0,1]$  范围内的值。即：

$$\tilde{h}_t = \tanh(C_t) \quad (2-18)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \quad (2-19)$$

其中， $W_{xo}$ 、 $W_{ho}$ 、 $b_o$  分别是输入向量、输出向量与输出门控向量的权重矩阵和偏置向量。

### (5) 输出 $h_t$

根据当前时间步的状态向量  $C_t$  和输出门  $o_t$ ，可以得到当前时间步的输出向量  $h_t$ 。具体地，我们将  $\tilde{h}_t$  乘上输出门  $o_t$  即可。即：

$$h_t = o_t \odot \tilde{h}_t \quad (2-20)$$

## (五) AT-LSTM 模型的基本理论

融合注意力机制的长短期记忆神经网络(Attention Long Short-Term Memory, AT-LSTM)是一种加入了注意力机制的时间序列预测模型，其核心思想是引入注意力机制来提高模型对于时间序列中不同时间步的关注程度，从而更好地捕捉时间序列的局部模式和全局趋势

[23]。

### 1. AT-LSTM 原理

AT-LSTM 的模型结构基于 LSTM 模型，在 LSTM 中，每个时间步的输入包括数据和前一个时间步的隐状态，经过门控单元的处理，得到输出和当前时间步的隐状态，从而实现对时间序列的建模。而在 AT-LSTM 中，为了更好地利用时间序列中不同时间步的信息，



引入了注意力机制。具体来说，将当前时间步的数据与前面隐状态进行融合，得到一个中间状态，然后通过一个注意力函数计算出每个时间步的权重，进而对所有时间步的中间状态进行加权平均得到最终的预测结果。这样可以使得模型更加关注重要的时间步，从而提高预测的准确性。AT-LSTM 神经网络结构图 2-2:

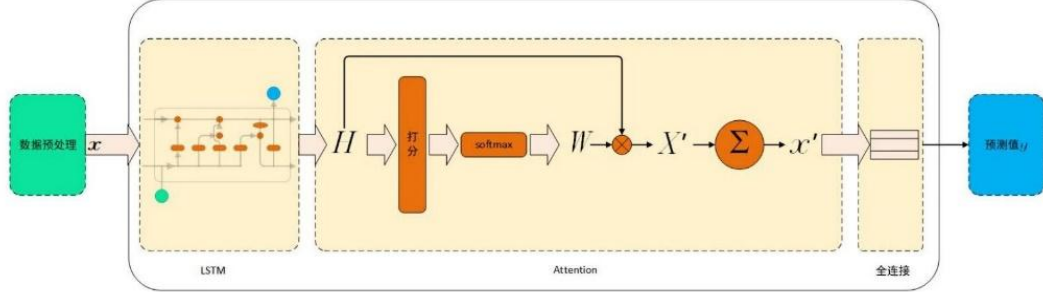


图 2-2 AT-LSTM 神经网络结构图

## 2. AT-LSTM 模型公式推导

本文对 AT-LSTM 模型中的公式进行推导，其中包括了 LSTM 的状态更新公式和注意力层的计算公式。

### (1) AT-LSTM 的状态更新公式

AT-LSTM 模型中的 LSTM 层与标准的 LSTM 模型类似，但是在状态更新公式中，需要将上下文向量  $r_t$  添加到输入向量  $[h_{t-1}, x_t]$  中。具体来说，LSTM 层的状态更新公式如下所示：

$$f_t = \sigma(W_f[h_{t-1}, x_t, r_t] + b_f) \quad (2-21)$$

$$i_t = \sigma(W_i[h_{t-1}, x_t, r_t] + b_i) \quad (2-22)$$

$$o_t = \sigma(W_o[h_{t-1}, x_t, r_t] + b_o) \quad (2-23)$$

$$\tilde{C}_t = \tanh(W_c[h_{t-1}, x_t, r_t] + b_c) \quad (2-24)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (2-25)$$

$$h_t = o_t \odot \tanh(C_t) \quad (2-26)$$

其中， $f_t$ 、 $i_t$ 、 $o_t$  分别是遗忘门、输入门和输出门的门控向量， $\tilde{C}_t$  是候选细胞状态， $C_t$  是细胞状态， $h_t$  是当前时间步的隐状态， $W_f$ 、 $W_i$ 、 $W_o$  和  $W_c$  是权重矩阵， $b_f$ 、 $b_i$ 、 $b_o$  和  $b_c$  是偏置向量， $\odot$  表示逐元素乘法。

### (2) 注意力机制的计算公式

注意力层的计算公式如下所示：

$$e_t = v^T \tanh(W_h h_{t-1} + W_x x_t + b_a) \quad (2-27)$$

$$\alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)} \quad (2-28)$$

$$r_t = \sum_{k=1}^T \alpha_k h_k \quad (2-29)$$

其中,  $e_t$  是每个时间步的注意力权重,  $v$  是用于计算注意力权重的权重向量,  $W_h$  和  $W_x$  是与 LSTM 层的隐状态和输入向量对应的权重矩阵,  $b_a$  是注意力层的偏置向量,  $\alpha_t$  是 Softmax 函数处理后的注意力权重,  $r_t$  是上下文向量, 通过加权平均计算得到。需要注意的是, 上述公式中的  $T$  是输入序列的长度,  $h_t$  是 LSTM 层在时间步  $t$  的隐状态,  $x_t$  是输入序列在时间步  $t$  的输入向量。

综上所述, AT-LSTM 模型的关键公式包括 LSTM 的状态更新公式和注意力层的计算公式。这些公式为 AT-LSTM 模型提供了对序列建模和上下文信息提取的能力, 使得模型可以更加准确地对序列进行分类、预测等任务。

### 三、数据集介绍

本文实验的数据集为通过雅虎财经网站获取的四支股票数据, 即中国平安股票(000001.SZ)、贵州茅台股票(600519.SS)、中国石油股票(601857.SS)和中国工商股票(601398.SS)。四个数据集选取的股票数据的时间跨度均为 2013 年 4 月 12 日至 2023 年 4 月 12 日、分别包含 2430 条股票数据。

其中中国工商股票数据集的部分数据如下表 3-1 所示:

表 3-1 中国工商股票数据集数据示例

Date	Open	High	Low	Close	Adj Close	Volume
2013/4/12	4.08	4.09	4.06	4.07	2.390317	33666531
2013/4/15	4.06	4.08	4.03	4.04	2.372698	43407375
2013/4/16	4.04	4.06	4.02	4.04	2.372698	39184563
2013/4/17	4.04	4.05	4.02	4.04	2.372698	38897888
2013/4/18	4.03	4.07	4.02	4.06	2.384444	29713214
2013/4/19	4.07	4.12	4.05	4.11	2.41381	70008889

四个股票数据集的数据字段含义如下表 3-2 所示:

表 3-2 股票数据各字段含义

字段	含义
Date	日期
Open	当日开盘价
High	当日最高价
Low	当日最低价
Close	当日收盘价
Adj Close	当日复权收盘价
Volume	当日成交量

本文实验任务是对四个股票数据集的收盘价进行预测，为了更直观的观察不同数据集收盘价的变化趋势，下面对四支股票数据的收盘价进行可视化，结果如图 3-1 至图 3-4 所示：



图 3-1 中国平安股票数据集收盘价变化趋势

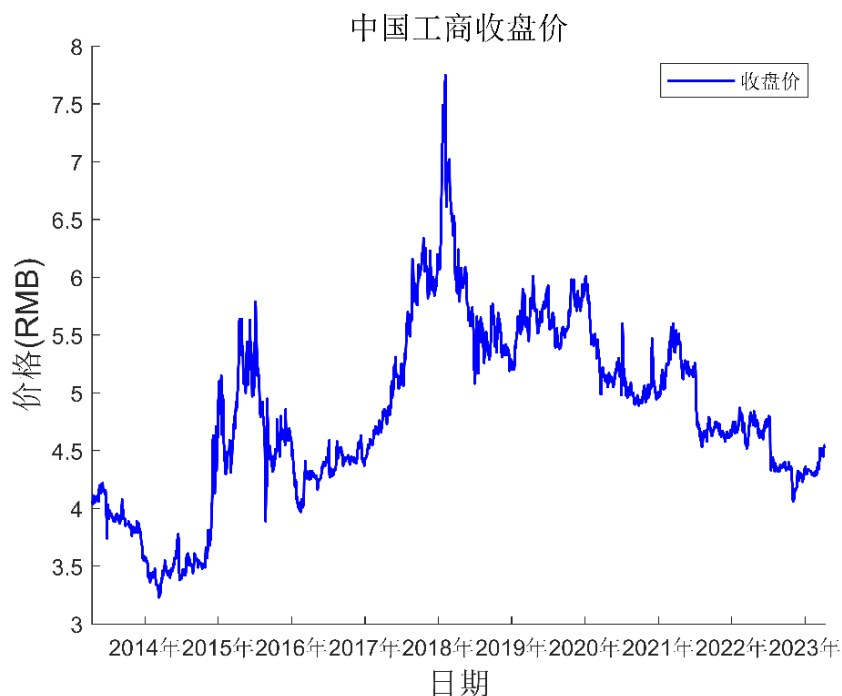


图 3-2 中国工商股票数据集收盘价变化趋势

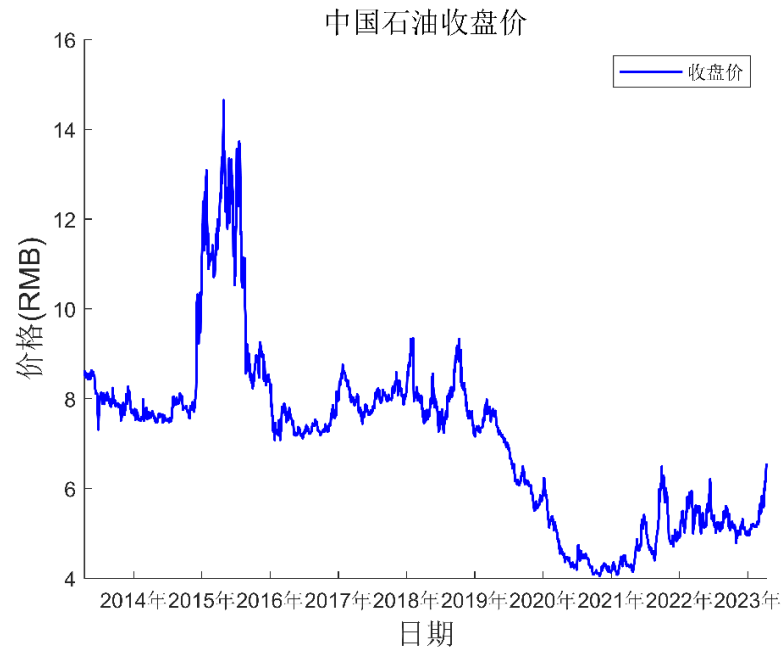


图 3-3 中国石油股票数据集收盘价变化趋势



图 3-4 贵州茅台股票数据集收盘价变化趋势

通过上述四支股票数据集收盘价变化趋势图，可以看出股票收盘价具有一定的周期性的，且在大周期下又存在着一定的波动。为了更直观的观察股票数据所具有的短期波动性，下面选取中国工商股票数据集 2023 年 1 月至 2023 年 3 月的部分数据进行可视化，结果如下图 3-5 所示：

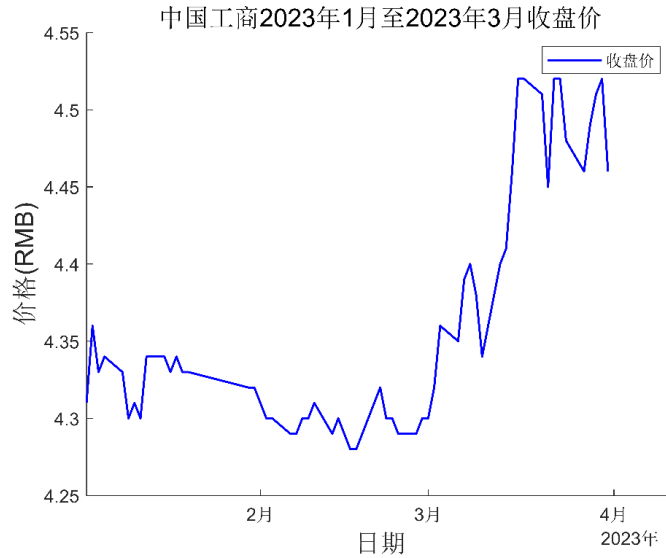


图 3-5：中国工商股票数据集短期收盘价变化趋势

为了更好的理解四支股票数据的特征和性质，以帮助我们做出更准确的分析和预测。首先计算出四个数据集收盘价的相关统计学指标，包括平均值、中位数、方差、标准差、偏度和峰度。各指标所表达的含义如下：

- (1) 平均值：通过股票收盘价的均值可以了解该股票在特定时间段内的平均表现水平。
- (2) 中位数：收盘价的中位数是指所有收盘价中位于中间的值，反映了股票的中间表现水平。由于中位数不受极端值的影响，因此其值比均值更稳定。
- (3) 方差：收盘价的方差是每个收盘价与其平均值的离差平方和的平均值，反映了股票收盘价的波动程度。方差越大，表明股票收盘价的波动越大。
- (4) 标准差：收盘价的标准差是方差的平方根，反映了股票收盘价的波动程度。标准差越大，表明股票收盘价的波动越大。
- (5) 偏度：偏度是一个统计量，反映了收盘价分布的不对称程度。正偏表示分布偏向较大值，负偏表示分布偏向较小值。
- (6) 峰度：峰度是一个统计量，反映了收盘价分布的峰度，即峰的陡峭程度。高峰度表示收盘价分布的峰较为陡峭，低峰度则表示峰相对平缓。

下面分别计算出四个股票数据集收盘价的六个统计学指标，结果如下表 3-3 所示：

表 3-3：收盘价统计学指标

数据集	平均值	中位数	方差	标准差	偏度	峰度
中国平安	12.05	11.20	17.50	4.18	<b>0.82</b>	3.27
中国工商	4.83	4.75	0.58	0.76	0.19	2.96
中国石油	7.10	7.53	3.82	1.95	0.71	<b>3.96</b>
贵州茅台	<b>856.46</b>	<b>665.10</b>	<b>478179.21</b>	<b>691.51</b>	0.55	1.77

从平均值和中位数指标上来看，贵州茅台股票数据的收盘价要明显比其他三支股票数据的收盘价高；从方差和标准差指标上来看，贵州茅台股票数据的收盘价波动性最强，中国工商股票数据的收盘价最为稳定；从偏度指标上来看，偏度值均为正数，说明四个数据集收盘价的数据分布均为正偏，即分布的尾部更长，数据更多地集中在左侧，且中国平安股票数据的收盘价的偏度值最高，说明其收盘价分布偏斜的程度比其他三个数据集大；从峰度指标上来看，峰度值均为正数，说明四个数据集收盘价数据分布的峰部比正态分布更尖锐，且中国石油股票数据的收盘价的峰度值最高，说明其收盘价分布峰态的陡峭程度比其他三个数据集更大。

## 四、基于 ARIMA 模型的金融数据变化趋势预测

### （一）数据预处理

#### 1. 数据平稳性检验

平稳性是指时间序列的统计特性不随时间变化而发生显著的变化，包括均值、方差和自相关性等指标，平稳性检验是用于检查一个时间序列是否具有平稳性质，一个平稳的时间序列具有稳定的统计规律，从而可以用较少的数据点进行建模和预测。由于传统的统计方法如回归分析、相关系数计算等，对于非平稳时间序列的可能不适用，因此，平稳性检验是时间序列分析的前提和基础，只有确保时间序列具有平稳性质，才能进行后续的分析、建模和预测。

常用的平稳性检验方法包括自相关函数和偏自相关函数的观察、单位根检验、差分运算等。本章实验选取单位根 ADF (Augmented Dickey-Fuller) 检验来判断数据是否平稳，可以检验序列是否存在单位根，即序列是否具有随机漫步性质。如果序列存在单位根，则说明序列非平稳；反之，如果序列不存在单位根，则说明序列是平稳的。下面对四个数据集进行单位根 ADF 检验以判断其平稳性，其中中国平安股票数据集收盘价平稳性检验过程如下：

滞后阶 (lag order) 用于描述一个时间序列中，一个观测值与其之前的观测值之间的滞后关系。在 ARIMA 模型中，滞后阶数指的是自回归项和移动平均项中的滞后阶数，本节实验使用 0-5 滞后阶对收盘价数据进行单位根 ADF 检验，中国平安股票数据集收盘价单根检验表结果如下表 4-1 所示：

表 4-1 中国平安股票数据集收盘价单根检验表

滞后阶	0	1	2	3	4	5
ADF 统计量	-0.23885	-0.24092	-0.23283	-0.25884	-0.26913	-0.23970
临界值	-1.9416	-1.9416	-1.9416	-1.9416	-1.9416	-1.9416
p 值	0.54447	0.56371	0.56667	0.55714	0.55338	0.56415

从 ADF 统计量上来看，0-5 滞后阶的 ADF 统计量的值都为负数，且绝对值都小于临界值。说明该数据序列具有单位根，即数据是非平稳的。从临界值上来看，0-5 滞后阶的临

界值均为-1.9416,这是由于在进行单位根检验时,使用的是常数项和趋势项均存在的假设。 $p$  值 ( $p$ -value) 反映了在原假设下,观察到的检验统计量的概率,如果  $p$  值小于设定的显著性水平 (通常为 0.05),则可以拒绝原假设,认为数据是平稳的,从结果可以看出,0-5 滞后阶的  $p$  值均大于 0.05,说明数据序列是非平稳的。

综上所述,中国平安股票数据集的收盘价数据是非平稳的,需要进行差分等操作以使其具有平稳性质。然后对另外三个数据集进行同样的单位根 ADF 检验,通过 ADF 统计量、临界值和  $p$  值综合分析其股票收盘价平稳性,结果如下表 4-2 所示:

表 4-2 股票收盘价平稳性表

股票名称	是否平稳
中国平安	否
中国工商	否
中国石油	否
贵州茅台	否

从上述实验结果可知,四支股票收盘价数据均为非平稳的时间序列,可能是由于序列中存在趋势或季节性等因素导致的,这些因素会导致序列的均值、方差或自相关性随时间发生变化,从而导致序列不具有平稳性质。由于本章实验使用的 ARIMA 模型要求数据为平稳性数据,因此需要对四个数据集进行预处理。

## 2. 数据平稳化

数据平稳化是将非平稳时间序列转化为平稳时间序列的一种常用方法。一般通过采用差分法、对数变换、移动平均等方法来处理非平稳时间序列,以使其具有平稳性质。

本节实验主要使用差分平稳化的方法对四个股票数据集的收盘价进行平稳化操作,其基本思想是对原始序列进行差分运算,从而使序列具有平稳性质。差分平稳化包括一阶差分、二阶差分等多种方法。其中,一阶差分是最常用的方法,它是将相邻两个数据点之间的差值作为新的数据点,得到一个新的序列。例如,对于一个非平稳化序列  $Y_1, Y_2, Y_3, \dots, Y_t$ , 其一阶差分序列为  $\Delta Y_1, \Delta Y_2, \Delta Y_3, \dots, \Delta Y_{t-1}$ , 其中  $\Delta Y_t = Y_t - Y_{t-1}$ 。下面对中国平安股票数据集的收盘价数据进行一阶差分,结果如图 4-2 所示。

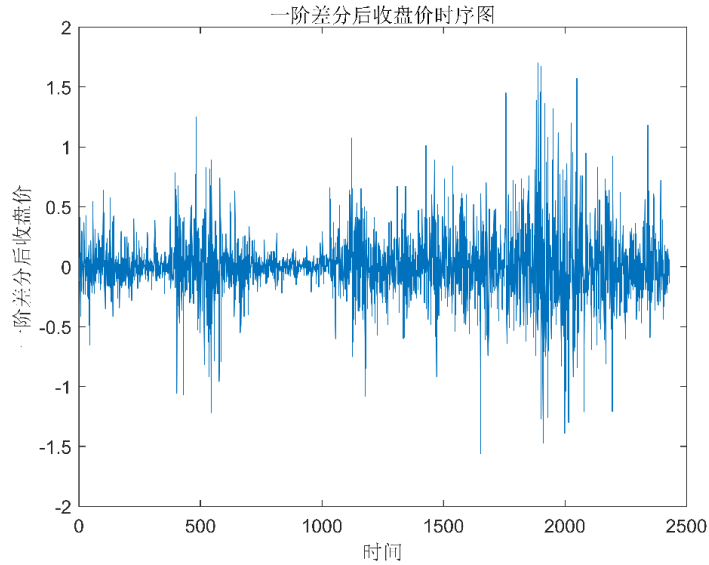


图 4-2 中国平安股票数据集收盘价一阶差分后时序图

然后对一阶差分后的数据单位根 ADF 检验，以判断其一阶差分后时间序列是否平稳，结果如下表 4-3 所示：

表 4-3 一阶差分后中国平安收盘价单根检验表

滞后阶	0	1	2	3	4	5
ADF 统计量	-48.8788	-35.115	-27.6132	-24.0256	-22.7458	-21.5753
临界值	-1.9416	-1.9416	-1.9416	-1.9416	-1.9416	-1.9416
p 值	0.001	0.001	0.001	0.001	0.001	0.001

上述实验结果显示，0-5 滞后阶的 ADF 统计量和 p 值都非常小，这意味着在所有滞后阶数下，ADF 统计量都拒绝了原假设，即序列不具有单位根，即一阶差分后的中国平安股票数据集收盘价数据序列是平稳的。

然后对其他三个数据集收盘价数据进行一阶差分，并对一阶差分后的时间序列进行单根 ADF 检验判断其平稳性，实验结果如下表 4-4 所示：

表 4-4 一阶差分后股票收盘价平稳性表

股票名称	是否平稳
中国平安	是
中国工商	是
中国石油	是
贵州茅台	是

综上所述，一阶差分后四支股票的收盘价时间序列数据都具有平稳性，因此可以用于后续基于 ARIMA 模型的预测。

## （二）实验环境及参数确定

### 1. 实验环境

本章实验的实验环境如下表 4-5 所示：



表 4-5 实验环境表

实验环境	详细情况
CPU	Intel(R) Core(TM) i7-9750H CPU @ 2.60GHz
GPU	NVIDIA GeForce GTX 1650 4G
内存	16.0 GB
系统	Windows11
软件	Python3.9+Jupyter+Matlab2018b

## 2. 参数确定

本节实验利用自相关函数（Autocorrelation Function, ACF）和偏自相关函数（Partial Autocorrelation Function, PACF）确定 ARIMA 模型的参数。ACF 表示时间序列中不同时间点之间的相关性。ACF 的值范围在 -1 到 1 之间，值越接近 1 表示两个时间点之间的相关性越强，值越接近 -1 则表示两个时间点之间的相关性越弱。PACF 表示在控制其它滞后项的情况下，两个时间点之间的相关性。PACF 的值范围在 -1 到 1 之间，值越接近 1 表示两个时间点之间的相关性越强，值越接近 -1 则表示两个时间点之间的相关性越弱。

其中对于 ARIMA(p,d,q)模型，其 ACF 和 PACF 图的特征如下：

- (1)自回归项 p: PACF 有显著的前 p 个峰，ACF 呈指数衰减。
- (2)移动平均项 q: ACF 有显著的前 q 个峰，PACF 呈指数衰减。
- (3)差分阶数 d: 差分一次或多次，直到得到平稳序列。

本节以中国平安股票收盘价为例，通过其 ACF 和 PACF 的函数图像确定预测中国平安股票收盘价 ARIMA 模型的参数。其 ACF 和 PACF 函数图像如下图 4-3 和 4-4：

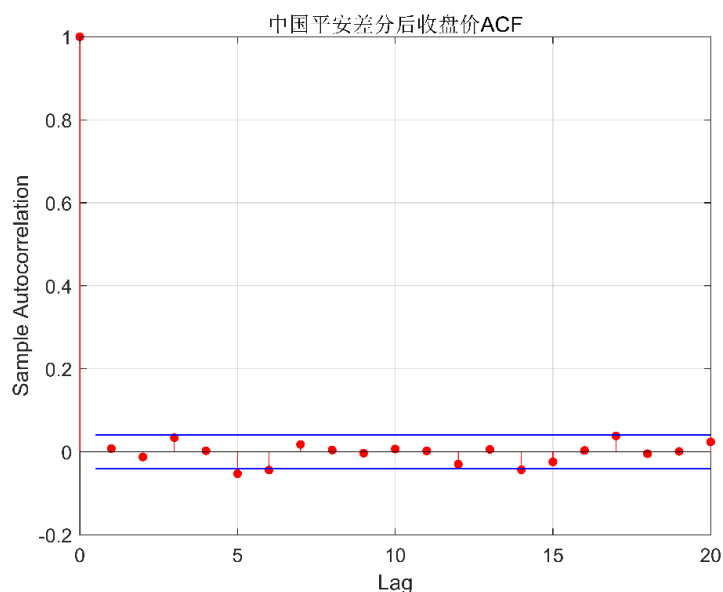


图 4-3：中国平安收盘价差分后 ACF 函数图

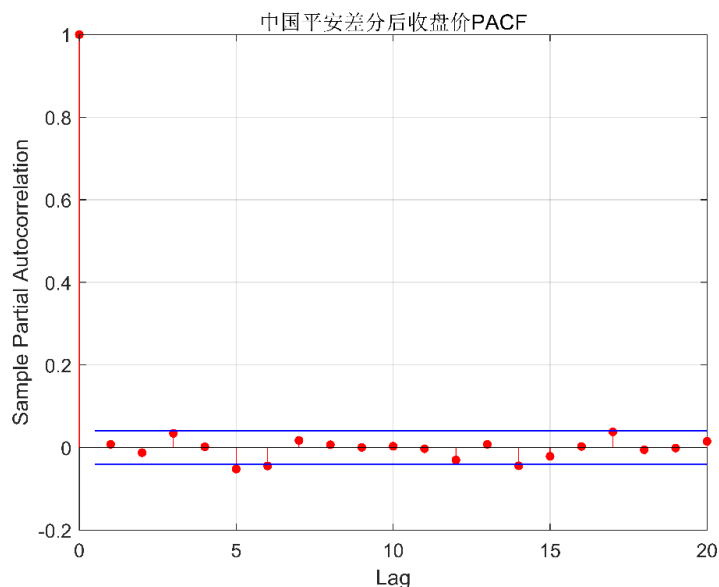


图 4-4: 中国平安收盘价差分后 PACF 函数图

通过上述函数图像，可以确定其 ARIMA 模型的最优参数组合为 ARIMA(1,1,1)。根据上述中国平安股票收盘价确定 ARIMA 模型步骤，可以确定其余三支股票 ARIMA 模型的最优参数组合为下表 4-6:

表 4-6: 各股票 ARIMA 模型最优参数组合表

股票名称	p	d	q
中国平安	1	1	1
中国工商	1	2	1
中国石油	1	1	1
贵州茅台	2	1	1

### （三）实验结果分析

分别对四支股票收盘价进行 ARIMA 模型预测实验，预测结果如图 4-5 至图 4-8 所示。

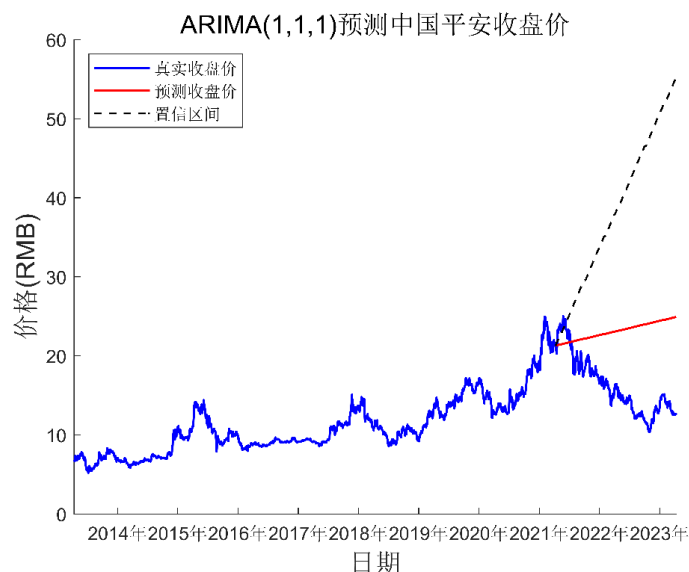


图 4-5 基于 ARIMA 预测中国平安股票收盘价

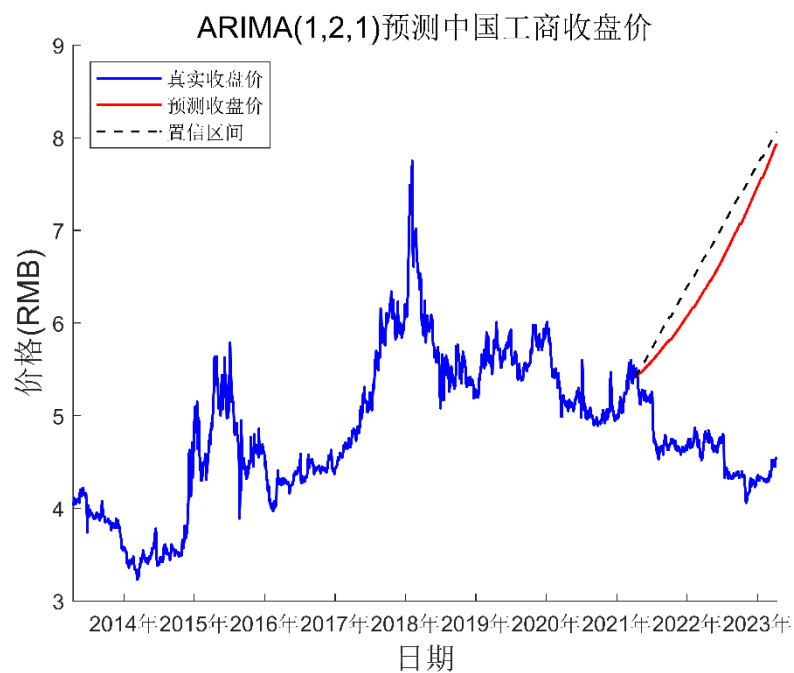


图 4-6 基于 ARIMA 预测中国工商股票收盘价

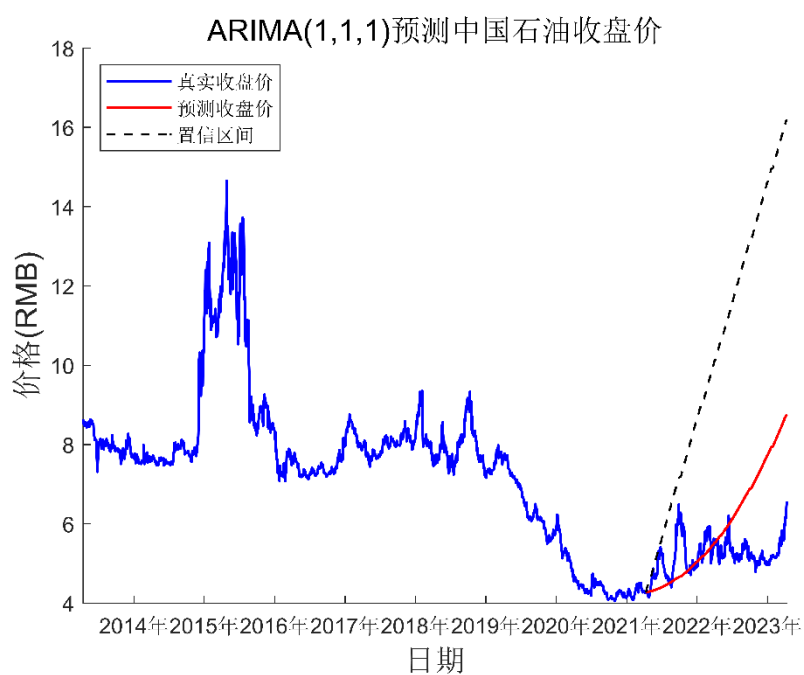


图 4-7 基于 ARIMA 预测中国石油股票收盘价

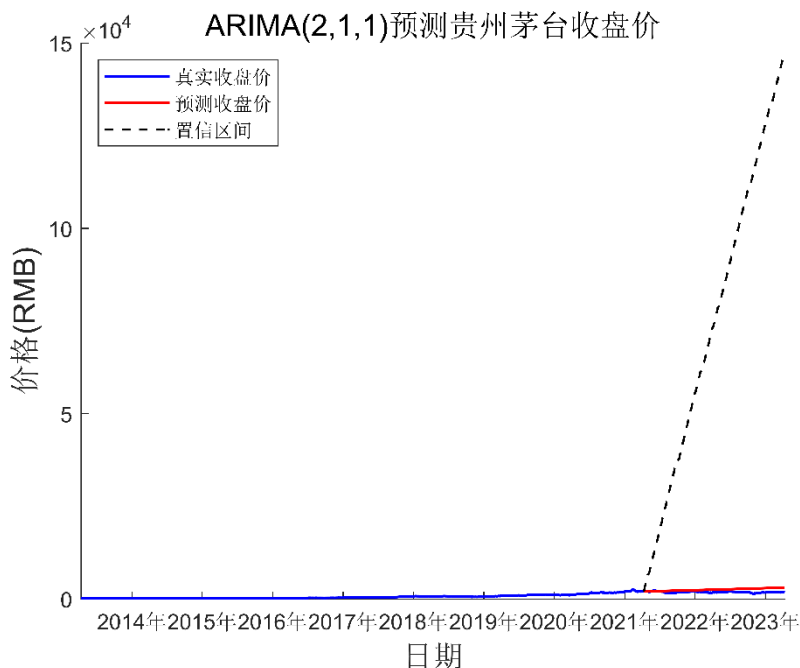


图 4-8 基于 ARIMA 预测贵州茅台股票收盘价

通过观察 ARIMA 模型预测各股票的对比图像，可以看出 ARIMA 模型在中国石油股票数据集上的预测效果较好，由于贵州茅台收盘价数据波动性大且周期性不明显，导致 ARIMA 模型在贵州茅台股票数据集上的预测效果较差。此外，从四个实验结果可以看出，ARIMA 模型在预测短期的股票收盘价上的预测效果比长期的好，因此要根据股票数据集具体的数据趋势、数据特征及预测时间序列的长短综合考虑选择合适的预测模型。

## 五、基于 LSTM 和 AT-LSTM 模型的金融数据变化趋势预测

### （一）数据预处理

本章基于 LSTM 和 AT-LSTM 模型在四个股票数据集上进行收盘价的预测实验，通过对数据集缺失值进行统计，发现四支股票运行状况良好，除节假日不开盘的情况之外，未有缺失数据，因此不必进行缺失数据的处理。此外，为了使模型能够更好地学习数据的模式，本章选择最大-最小归一化（Min-Max Normalization）方法来处理数据集。最大-最小归一化也称为离差标准化，是一种常见的数据归一化方法，能够将数据按照一定比例缩放到指定的区间内。该方法的原理是将数据的最小值设为 0，最大值设为 1，然后按照线性比例将数据缩放到[0,1]之间。通过最大-最小归一化操作后的数据，可以消除量纲影响、提高模型收敛速度、减小异常值的影响和提高数据可视化效果。

下面首先对四个股票数据集进行最大-最小归一化处理，其中中国平安股票数据集收盘价数据经过最大-最小归一化处理后的结果如下图 5-1 所示：



图 5-1 中国平安收盘价归一化图

## (二) 基于 LSTM 模型的实证分析

### 1. 实验环境

本章实验的实验环境如下表 5-1 所示：

表 5-1 实验环境

实验环境	详细情况
CPU	Intel(R) Core(TM) i7-9750H CPU @ 2.60GHz
GPU	NVIDIA GeForce GTX 1650 4G
内存	16.0 GB
系统	Windows11
软件	Python3.9+pandas1.4.2+pytorch1.11.0+Jupyter

### 2. 模型搭建

#### (1) 参数介绍

模型各参数的含义如下表 5-2 所示：

表 5-2 LSTM 模型参数含义表

参数名称	参数设定值	含义
input_dim	1	数据输入特征数量
hidden_dim	32	每个隐藏层的神经元数量
num_layers	2	隐藏层层数
output_dim	1	预测输出特征数量
num_epochs	200	训练轮数
lr	0.01	学习率

## (2)模型实现

本节实验基于 PyTorch 框架，首先对数据进行最大-最小归一化处理，得到数据后对数据进行适应化改造，以适合 LSTM 模型输入格式的数据。划分出总数据的 80%为训练集，20%为测试集，然后设置模型的参数，选择 MSE 作为损失函数以对模型进行优化。

## (3)模型搭建过程

**准备数据集：**在开始训练 LSTM 之前，需要准备一个适当的数据集。这通常涉及到数据的清洗、预处理和划分成训练集、验证集和测试集。

**构建模型：**在准备好数据集之后，就可以开始构建 LSTM 模型。通常可以使用深度学习框架如 PyTorch 来实现。在构建模型时，需要选择适当的网络结构和超参数。

**编译模型：**在构建完 LSTM 模型之后，需要编译它以便可以开始训练。编译模型涉及到选择适当的损失函数、优化器和指标。

**训练模型：**训练 LSTM 模型通常涉及到迭代地将数据送入模型中，并根据训练数据的反馈来更新模型参数。在训练期间，需要监测模型在验证集上的性能，以便可以进行超参数调整和早期停止。

**评估模型：**在训练 LSTM 模型之后，需要对其进行评估以检查其在测试集上的性能。这通常涉及到计算一些指标，如 MSE、MAE、RMSE。

**使用模型：**在评估完 LSTM 模型之后，可以将其用于实际预测。这涉及到将新数据输入到模型中，并使用训练得到的参数来生成预测结果。

## 3. 实验结果分析

本节基于 LSTM 模型在四个股票数据集上进行收盘价的预测实验，使用 MSE、RMSE 和 MAE 三个模型评价指标来评价模型预测股票收盘价的准确度。下图 5-2 是基于 LSTM 的中国平安股票数据集收盘价预测结果的损失函数 MSE 变化图。

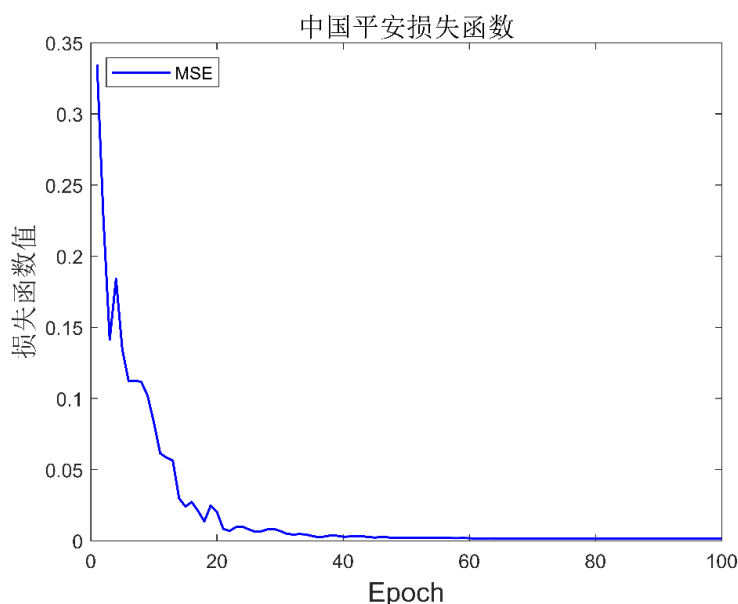


图 5-2 基于 LSTM 的中国平安股票收盘价预测损失函数变化图

从上图可以看出损失函数不断趋向 0，表示模型对训练集的拟合越来越好，最终经过 100 轮的训练，MSE 的值趋向于 0，考虑到模型的损失函数趋向于 0 并不代表着预测效果一定好，因此使用训练好的 LSTM 模型对后 20% 的测试集数据进行预测，结果如下图 5-4 所示：



图 5-3 基于 LSTM 的中国平安股票收盘价预测结果图

通过实验得出其余三支股票数据集收盘价预测结果的损失函数 MAE、MSE 和 RMSE，结果如下表 5-3 所示：

表 5-3 基于 LSTM 模型的评价指标结果

股票名称	MAE	MSE	RMSE
中国平安	0.4042	0.3053	0.5225
贵州茅台	46.9256	3601.222	60.0102
中国石油	0.1212	0.0275	0.1659
中国工商	0.0565	0.0095	0.0973

从结果可以看出，中国平安和中国石油的 MAE、MSE 和 RMSE 均较小，说明预测结果比较准确；而贵州茅台的 MAE、MSE 和 RMSE 均较大，说明预测误差较大。中国工商的 MAE 和 MSE 较小，但 RMSE 较大，说明预测值与真实值之间的差异比较大，需要进行进一步的优化。

### （三）基于 AT-LSTM 模型的实证分析

#### 1. 模型搭建

##### （1）参数介绍

模型各参数的含义如下表 5-4 所示：

表 5-4 AT-LSTM 模型参数含义表

参数名称	参数设定值	含义
input_dim	1	数据输入特征数量
hidden_dim	32	每个隐藏层的神经元数量
num_layers	2	隐藏层层数
output_dim	1	预测输出特征数量
num_epochs	100	训练轮数
batch_size	64	批次大小
lr	0.01	学习率

## (2)模型实现

本节实验同样基于 PyTorch 框架，首先对数据进行最大-最小归一化处理，得到数据后对数据进行适应化改造，以适合 LSTM 模型输入格式的数据。划分出总数据的 80%为训练集，20%为测试集，然后设置模型的参数，同时使用 MSE 当作损失函数，对模型进行优化。

## 2. 实验结果分析

本节基于 AT-LSTM 模型在四个股票数据集上进行收盘价的预测实验，使用 MSE、RMSE 和 MAE 三个模型评价指标来评价模型预测股票收盘价的准确度。下图 5-4 是基于 AT-LSTM 的中国平安股票数据集收盘价预测结果的损失函数 MSE 变化图。

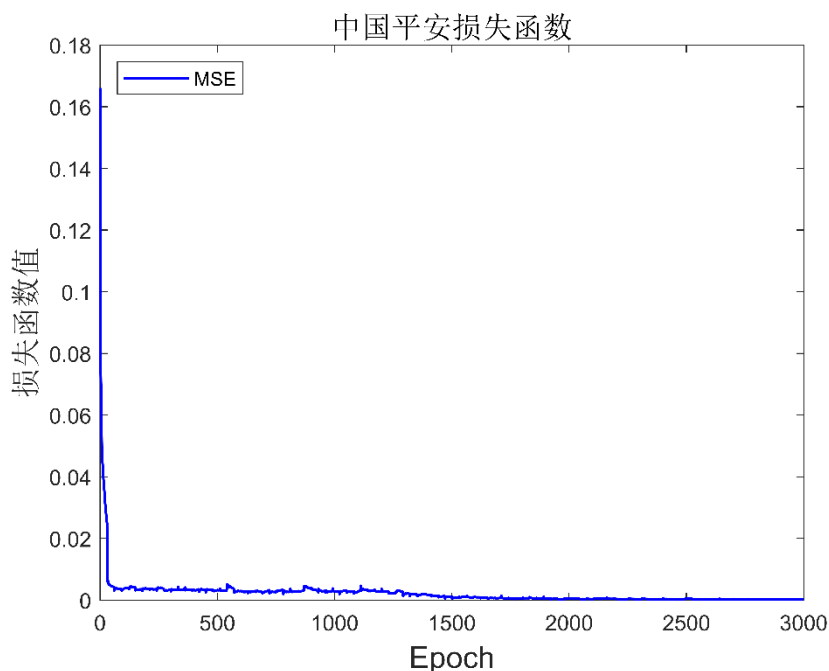


图 5-4 基于 AT-LSTM 的中国平安股票收盘价预测损失函数变化图

从上图可以看出损失函数不断趋向 0，表示模型对训练集的拟合越来越好，最终经过 100 轮的训练，MSE 的值趋向于 0，考虑到模型的损失函数趋向于 0 并不代表着预测效果一定好，因此使用训练好的 LSTM 模型对后 20%的测试集数据进行预测，结果如下图 5-5 所示：



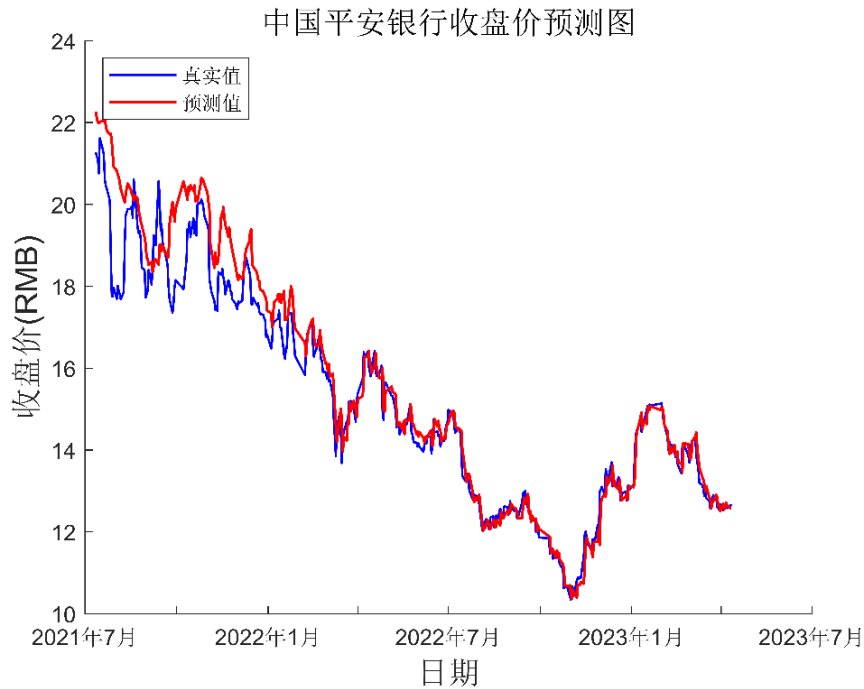


图 5-5 基于 AT-LSTM 的中国平安股票收盘价预测结果图

通过实验得出其余三支股票预测后计算的 MAE、MSE 和 RMSE 值，结果如下表 5-5 所示：

表 5-5 基于 AT-LSTM 模型的评价指标结果

股票名称	MAE	MSE	RMSE
中国平安	0.4735	0.5635	0.7507
贵州茅台	82.8673	12020.092	109.6362
中国石油	0.2854	0.1506	0.3881
中国工商	0.0558	0.0060	0.0775

对于中国平安、中国石油和中国工商三支股票数据集来说，它们的 MAE、MSE 和 RMSE 值都比较小，说明预测效果较好。其中，中国工商的预测效果最好，MAE、MSE 和 RMSE 均最小，分别为 0.4735、0.5635、0.7507。对于贵州茅台股票数据集来说，其 MAE、MSE 和 RMSE 值都比较大，特别是 RMSE 值高达 109.6362，说明预测效果较差，可能是由于贵州茅台股票价格波动较大，存在较多异常值，使得预测模型的效果受到影响。

综上所述，使用 ARIMA、LSTM、AT-LSTM 三种模型在中国平安股票、中国石油股票和中国工商股票三个数据集上都能够较好的完成对收盘价进行预测的任务，且 AT-LSTM 模型的预测效果最好。但由于贵州茅台股票数据的波动性较大、周期性不明显，导致三种模型在该数据集上的预测效果较差，因此在金融数据预测任务上，要根据股票数据集具体的数据趋势、数据特征及预测时间序列的长短综合考虑选择合适的预测模型，以达到更好的预测效果。

## 六、总结与展望

### （一）研究总结

本文首先对金融时间序列预测进行了介绍，并介绍了序列数据的特点和性质，同时对国内外金融时间序列数据预测的研究现状进行阐述，针对预测存在的问题和可能性进行研究。然后介绍了金融时间序列数据常用的预测方法和模型，对 ARIMA、LSTM、AT-LSTM 三个模型的结构及原理进行了详细介绍，并在中国平安股票、中国石油股票、中国工商银行和贵州茅台股票四个股票数据集上进行收盘价预测的实验，通过实验结果进行比较分析，实验结果显示，AT-LSTM 模型在中国平安股票、中国石油股票和中国工商银行三种股票数据集上的预测效果最好，但由于贵州茅台股票数据的波动性较大、周期性不明显，导致 ARIMA、LSTM、AT-LSTM 三种模型在该数据集上的预测效果都较差。因此，在金融数据预测任务上，要根据股票数据集具体的数据趋势、数据特征及预测时间序列的长短综合考虑选择合适的预测模型，以达到更好的预测效果。

本文的主要工作总结如下：

1.建立了基于 ARIMA 金融数据变化趋势的预测模型。通过利用 ARIMA 模型的原理分析，并对实验数据进行分析，利用 ARIMA 模型的原理，对数据进行平稳性检验，再通过一阶差分，是数据平稳，再通过 ARIMA 模型的预测功能实现与时间序列数据的预测。

2.建立了基于 LSTM、AT-LSTM 神经网络的金融数据变化趋势的预测模型。并分析了两个模型的原理和区别，着重介绍了 AT-LSTM 模型引入注意力机制，并对模型进行实验，得出了 AT-LSTM 模型比 LSTM 预测金融时间序列数据更加精准的结论。

3.通过对比分析 ARIMA、LSTM、AT-LSTM 模型的实验结果，对 ARIMA 模型、LSTM 模型和 AT-LSTM 模型进行了综合评价，并提出了未来时间序列模型研究的方向。ARIMA 模型具有简单易用和强解释性等优点，但是对于非线性和非稳定的数据表现不佳。LSTM 模型能够处理非线性和长期依赖的数据，但是在噪声和非平稳的数据上表现欠佳。相比之下，AT-LSTM 模型更具有鲁棒性和可解释性，但是需要更长的训练时间和更大的计算量。

### （二）研究展望

1.数据特征较少，只有中国平安股票的收盘价、开盘价、最高价、最低价、交易量和经调整的收盘价。因此未来可以引入其他的特征，比如该行业的板块效应，新闻报道的文本情感，大盘指数等特征。

2.本文在对金融时间序列数据预测的时候过于理想化，因此我们需要增加实际因素，增大时间序列的长度以确保股票价格发展时间的完整性。比如 2008 年金融危机等状况。

3.本文利用 ARIMA 模型、LSTM 模型、AT-LSTM 模型实现对股票价格预测，但是数据量，特征选择较少，同时可能出现过拟合的现象。因此以后在时间序列预测的神经网络模型结构中预测时需要考虑过拟合情况。

4.时间序列数据预测是机器学习和统计分析中的一个重要领域，它涉及使用历史数据

来预测未来的趋势、周期性变化和异常值。这个领域的应用广泛，包括金融、交通、天气预报、股票市场、能源消耗和销售预测等。需要我们不断探索新的技术和方法来提高预测精度和效率。

## 参考文献

- [1] Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
- [2] Tay, Francis EH, and L. J. Cao. "Modified support vector machines in financial time series forecasting." *Neurocomputing* 48.1-4 (2002): 847-861.
- [3] Hassan, Md Rafiul, et al. "A hybrid of multiobjective Evolutionary Algorithm and HMM-Fuzzy model for time series prediction." *Neurocomputing* 81 (2012): 1-11.
- [4] Zahedi, Javad, and Mohammad Mahdi Rounaghi. "Application of artificial neural network models and principal component analysis method in predicting stock prices on Tehran Stock Exchange." *Physica A: Statistical Mechanics and its Applications* 438 (2015): 178-187.
- [5] Lahmiri, Salim. "Minute-ahead stock price forecasting based on singular spectrum analysis and support vector regression." *Applied Mathematics and Computation* 320 (2018): 444-451.
- [6] Kim, Taewook, and Ha Young Kim. "Forecasting stock prices with a feature fusion LSTM-CNN model using different representations of the same data." *PloS one* 14.2 (2019): e0212320.
- [7] Luo, Shihua, and Cong Tian. "Financial high-frequency time series forecasting based on sub-step grid search long short-term memory network." *IEEE Access* 8 (2020): 203183-203189.
- [8] 裴双喜. 基于数据挖掘的金融时间序列预测分析与研究[D]. 大连海事大学, 2008.
- [9] 张文霄. 基于 PSO 优化的 BP 神经网络股票预测模型[D]. 哈尔滨工业大学, 2010.
- [10] 孙瑞奇. 基于 LSTM 神经网络的美股股指价格趋势预测模型的研究[D]. 首都经济贸易大学, 2016.
- [11] 王谨平. 基于金融数据的时间序列研究与应用[D]. 电子科技大学, 2017.
- [12] 周凌寒. 基于 LSTM 和投资者情绪的股票行情预测研究[D]. 华中师范大学, 2018.
- [13] 陈璐. 基于 LSTM 模型的金融时间序列预测算法研究[D]. 哈尔滨工业大学, 2019.
- [14] 田晓丹. 基于 LSTM 与多 GARCH 型混合模型的股价波动性预测的实证分析[D]. 哈尔滨工业大学, 2019.
- [15] 赵薇. 基于 LSTM 神经网络的金融数据预测分析[D]. 哈尔滨工业大学, 2020.
- [16] 李庆涛. 基于板块效应以及多元特征的金融时间序列预测研究[D]. 山东财经大学, 2022.
- [17] 于明加. 市场有效性理论研究综述[J]. *辽宁经济*, 2020, No. 435(06): 56-57.
- [18] 闫洪举. 金融时间序列数据预测: 文献回顾与展望[J]. *金融教育研究*, 2021, 34(03): 33-41.
- [19] 张倩玉. 神经网络在股票价格预测模型的研究与应用[D]. 天津财经大学, 2021.
- [20] 文海涛, 倪晓萍. 我国上市公司财务指标与股价相关性实证分析[J]. *数量经济技术经济研究*, 2003(11): 118-122.
- [21] 熊政, 车文刚. ARIMA-GARCH-M 模型在短期股票预测中的应用[J]. *陕西理工大学学报(自然科学版)*, 2022, 38(04): 69-74.
- [22] 武博. 基于 LSTM 模型的股票价格预测[D]. 大连理工大学, 2021.
- [23] 张怡. 基于 ARIMA 和 AT-LSTM 组合模型的股票价格预测[J]. *电脑知识与技术*, 2022, 18(11): 118-121.

## 致谢

中国山东省济南市天桥区大桥镇司家村，司佳小学。

中国山东省济南市天桥区大桥镇靳家村，大桥镇第二中学。

中国山东省济南市历城区郭店镇三区 20 号，山东省实验中学。

中国山东省济南市历下区二环东路 7366 号，山东财经大学。

Sverige Uppsala 751 05 Uppsala Box 256, Uppsala Universitet.