

项目编号: IPPXQ01004

上海交通大学

大学生创新创业计划

项目研究论文

论文题目: 量化对冲基金的高频交易技术方法论研究

项目负责人: 龚逸炜 学院(系): 数学科学学院

指导教师: 任桐鑫 学院(系): 学生创新中心

参与学生: 龚逸炜, 韩圣依, 殷奕珉, 唐华

项目执行时间: 2022 年 3 月 至 2023 年 5 月

摘要

高频交易 (High Frequency Trading) 是一种新型金融交易模式, 使用高频率的算法交易进行实时分析和快速交易, 近年来发展迅速。本研究旨在探讨高频交易技术在量化对冲基金中的应用和运作机制。通过理论研究和实证分析, 深入研究高频交易技术的优缺点, 并设计适用于高频交易的模型, 根据市场数据对比评估, 为制定更有效的投资策略和风险管理措施提供参考。

本研究采用了两种典型的模型进行预测和对比分析。

第一种模型是 LSTM 模型, 它是一种循环神经网络模型, 可以处理非线性关系和非平稳时间序列。在本研究中, LSTM 模型将多种市场因子纳入考虑, 并将整体预测任务分为两部分: 股价具体值预测部分和股价涨跌符号预测部分。通过超参数寻优和网络训练, 结果显示, 在股价涨跌符号预测任务上 LSTM 模型的表现较好, 测试集上预测准确度较高, 趋势预测任务比值预测任务更容易。

第二种模型是 ARIMA 模型, 它是一种线性模型, 需要较少的参数来建立模型, 适用于稳定的时间序列数据。本研究使用了一种新颖的策略, 将 $t - k_1$ 到 t 时刻的收益作为动量因子、 $t - k_2$ 到 t 时刻的随机波动作为反转因子, 以此来预测 t 到 $t + 2$ 时刻的收益。这种组合方法能够更全面地考虑股票市场的复杂性, 权衡精度和简便性。

根据模型对比分析得到: LSTM 模型适用于高频交易的股价涨跌预测任务, 对于较为稳定的时间序列数据具有较好的预测效果, 但在处理噪声较大的非平稳时间序列数据时, 可能会出现过拟合的问题。ARIMA 模型能够更好地提取序列数据中的趋势和周期性, 适用于稳定的时间序列数据, 但在处理非线性关系的序列数据时性能较差, 可能会出现欠拟合的问题。总体来说, LSTM 模型和 ARIMA 模型都有其适用范围和优缺点, 需要根据实际情况选择合适的模型。

本研究通过理论研究和实证分析, 深入探讨了高频交易技术在量化对冲基金中的应用和运作机制, 比较了 LSTM 模型和 ARIMA 模型在高频交易中的优缺点和适用范围。通过对两种模型的分析 and 对比, 得出结论: 在高频交易中, LSTM 模型和 ARIMA 模型可以结合使用, 提高预测精度和可靠性。

随着高频交易技术和量化对冲基金的不断发展, 我们可以进一步探索更高效、准确的预测模型, 改善模型解释性和数据安全性等问题, 优化投资策略和风险管理。

关键词: 高频交易, 机器学习, LSTM, ARIMA

ABSTRACT

High Frequency Trading (HFT) is a new financial trading mode, which has developed rapidly in recent years. This study aims to explore the application and operation mechanism of HFT technology. Through theoretical research and empirical analysis, we investigate advantages and disadvantages of HFT technology and design models suitable for HFT. By comparing and evaluating real data, this study provides a reference for investment strategies.

Two typical models were used for prediction and comparative analysis in this study.

The first model is LSTM model, a type of recurrent neural network model that can handle nonlinear relationships and non-stationary time series. We consider multiple market factors to both predict the specific value of stock prices and the direction of stock price movements. Results show it performs better in predicting the direction of stock price movements.

The second model is ARIMA model, a linear model that is suitable for stable time series data. This study used a novel strategy, taking the return from $t - k_1$ to t as the momentum factor and the random fluctuations from $t - k_2$ to t as reversal factor to predict the returns from t to $t + 2$. This combination method comprehensively balance accuracy and simplicity.

According to comparative analysis of the models, the LSTM model is suitable for predicting the direction of stock price movements while it may overfit when dealing with non-stationary time series data. The ARIMA model better extracts trends in stable sequence data while its performance is poor when dealing with non-linear relationships in sequence data. Both the LSTM model and ARIMA model have their applicable ranges, and the appropriate model should be selected based on the actual situation.

Through theoretical research and empirical analysis, this study thoroughly explores application and mechanism of HFT and compares the advantages and disadvantages of LSTM model and ARIMA model. By comparing two models, the conclusion is drawn that these two models can be combined in HFT to improve prediction accuracy and reliability.

As HFT technology and quantitative hedge funds continue to develop, we can further explore more efficient and accurate prediction models, improve model interpretability and data security, and optimize investment strategies.

Key words: High Frequency Trading, Machine Learning, LSTM, ARIMA

目 录

第一章 绪论	1
1.1 研究背景	1
1.2 研究意义	1
1.3 文献综述	2
1.4 研究内容	3
1.5 论文创新点	4
第二章 理论模型和研究方法	5
2.1 高频交易概述	5
2.1.1 交易特点	5
2.1.2 发展背景	5
2.1.3 策略研究	6
2.1.4 交易工具	7
2.2 长短期记忆神经网络 (LSTM)	7
2.2.1 循环神经网络 (RNN)	8
2.2.2 LSTM 基本结构	8
2.2.3 LSTM 工作原理	9
2.3 ARIMA 模型.....	10
2.3.1 自回归模型 AR.....	11
2.3.2 移动平均模型 MA	11
2.3.3 差分处理 I	11
2.3.4 ARIMA 模型的建立	11
2.3.5 模型优缺点	12
第三章 量化策略与实证分析	14
3.1 数据集介绍	14
3.2 模型假设	14
3.3 模型构建	15
3.3.1 LSTM 模型.....	15

3.3.2	ARIMA 模型.....	18
3.4	数据分析	20
第四章	研究结论与展望	23
4.1	研究结论	23
4.1.1	LSTM 模型.....	23
4.1.2	ARIMA 模型.....	23
4.1.3	模型对比	24
4.2	不足与展望	24

第一章 绪论

1.1 研究背景

金融市场是比较有效的市场，但其中仍然存在大量的非理性定价空间，优秀的量化高频交易策略和技术有助于提升金融交易市场的有效性、合理性。高频交易（High Frequency Trading）是近年来出现的一种创新金融交易模式。欧洲证券监管委员会对高频交易的定义是：利用复杂的计算机与 IT 系统，以毫秒级的速度执行交易并日内短暂持有仓位，以超高的速度在不同的交易平台进行交易的自动化交易形式。在外汇市场中，高频交易是指以高周期率和高订单率的算法交易（Algorithmic Trading），对市场数据进行实时分析，在市场发生变化时迅速执行相应的交易策略。美国商品期货交易委员会将高频交易定义为一种通过预先设定的程序算法实现的高速度、高频次的交易模式。随着计算机技术和数据分析方法的不断发展，高频交易技术不断演化和升级，为量化对冲基金带来了更高的收益率和更低的风险。因此，对于量化对冲基金的高频交易技术方法论的研究具有重要的理论和实践意义。

近几年来，高频交易（high frequency trading）在全球金融市场得到了迅速发展。这高频交易作为电子化交易的最新模式，其高盈利性与高争议性的特点引发了业界与学术界的广泛关注。高频交易的核心技术之一，便是用于预测未来短期内市场价格的变化。这是一门技术性很强的学科，它吸引了来自数学、物理、计算机科学、电子工程等领域中的顶尖人才，号称交易领域的“金字塔尖”。从目前高频交易策略的研究现状来看，通常有三类策略，其中两类是服务于交易客户和为市场提供流动性的巨量订单拆解策略和做市流动性交易策略，还有一类是关注盈利性的定量交易策略。然而，随着行业的发展与算法的迭代，高频交易的运营成本提高，利润潜力下降，故而当下行业最重要的问题是如何开发更优化的算法。高频交易市场正在不断完善过程中，学界和业界也通过不同模型进行假设分析，值得我们进一步的学习研究和实践分析。

1.2 研究意义

本研究旨在探讨高频交易技术在量化对冲基金中的应用，分析模型应用的优缺点及其适用的范围。通过对量化对冲基金的高频交易技术进行深入研究，可以更好地理解量化对冲基金的运作机制，为量化对冲基金投资者提供更多的信息，帮助其更

好地理解和把握市场风险。同时，研究结果对于投资机构和金融机构存在参考价值，可以帮助其制定更有效的投资策略和风险管理措施。

1.3 文献综述

量化交易在国外已经发展 40 多年，已趋于成熟，交易品种相对多样交易市场完善。机器学习算法在国外量化市场上的应用已经很广泛。国外相关研究主要集中在机器学习在股票市场中的应用，少部分聚焦于期货市场研究。研究早期的算法主要以鲁棒性较好的支持向量机、随机森林、梯度提升树等算法为主。近些年来，伴随着深度学习的发展，有关 LSTM 算法的研究越来越多，但也基本集中于股票市场。

德国的 Gen Cay 在 1996 年发表的《Non-linear prediction of security returns with moving average rules》，使用前馈人工神经网络模型，基于 1967 到 1988 年的道·琼斯工业指数历史数据的 7 日均线数据上，进行了神经网络在金融数据上的应用研究。实验结果证明，前馈人工神经网络模型的预测精度比传统统计学的移动平均算法更高。M. Karazmodeh, S. Nasiri 和 S. Majid Hashemi 在 2013 年发表的《Stock Price Forecasting using Support Vector Machines and Improved Particle Swarm Optimization》中，基于 SVM 的框架，使用遗传算法（IPSO）改良的粒子群优化算法（PSO），对各种股票的指数进行有效预测，计算效率和预测精度较高，具有较强的鲁棒性。

David M. Q. Nelson, Adriano C. M. Pereira 等在 2017 年发表的《Stock market's price movement prediction with LSTM neural network》中运用 LSTM 算法对股票进行趋势预测。以历史价格数据为基础，结合技术分析指标，运用 LSTM 网络预测股票在未来一段时间内的走势。实验结果表明，这一模型在预测特定股票的涨跌情况的准确率平均能达到 55.9%。

Sreelekshmy Selvin, Vinayakurnar R 等在 2017 年发表的《Stock price prediction using LSTM, RNN and CNN-sliding window model》中指出，股票市场对经济存在深远的影响。现有的预测方法利用了线性算法（AR、MA、ARIMA）和非线性算法（ARCH、GARCH、神经网络），但这些算法侧重于利用每日收盘价来预测单个公司的股票指数走势。而该研究中，采用深度学习体系结构来识别数据中的潜在动态，应用三种不同的深度学习架构，实现股价预测。

孙瑞奇在 2015 年发表的《基于 LSTM 神经网络的美股股指价格趋势预测模型的研究》中，研究分析前馈神经网络、循环神经网络和 LSTM 神经网络对股票短期价格作预测的可行性并进行了对比和改良。实验结果证明，LSTM 神经网络模型可以对股

票历史数据进行学习,找出历史时间序列数据之间的非线性关系,学习出某些影响价格走势的特征因子,并深度挖掘股票价格时间序列中的固有规律,实现了股票短期价格预测。

王苏生等在 2018 年发表的《基于 ARMA 模型的沪深 300 股指期货高频数据收益率研究与预测》中,选择了沪深 300 股指期货 10 天内每 0.5 秒一次的高频数据作为研究对象,分析了高频数据日内收益率分布特征,发现高频数据在日内有波动率聚集现象,但分布并没有出现尖峰现象和后尾现象。文章通过进行自相关分析与偏自相关分析,发现日内高频数据收益率存在较强的自相关性。论文采用 EACF 方法对训练样本定阶,确定了 ARMA 模型的参数,并通过最小二乘法计算出了各模型系数。实验结果表明,ARMA 模型能够较好地模拟沪深 300 股指期货日内收益率的曲线走势。

黄卿、谢合亮在 2018 年发表的《机器学习方法在股指期货预测中的应用研究——基于 BP 神经网络、SVM 和 XGBoost 的比较分析》中,针对股指期货交易速度快、交易频率高和交易量巨大、时序性等特征,构建了新的股指期货量化投资模型,采用沪深 300 股指期货 1 分钟粒度数据作为研究对象。并对比分析了神经网络、支持向量机和 XGBoost 对股指期货下 1 分钟价格的变动方向的预测能力。研究结果证明,三种机器学习方法都具有较好的预测能力,其中 XGBoost 的预测能力要优于传统的神经网络和 SVM 模型。

纵观国内外的文献研究可以发现,金融数据序列预测数据早年的研究算法以线性预测和传统的机器学习为主。传统的线性预测方法,如自回归模型 (AR)、滑动平均模型 (MA)、自回归滑动平均模型 (ARMA) 都是基于历史序列数据、白噪声或两者组合的线性预测方法,但要求时间序列是平稳的或差分序列是平稳的,因此存在一定的局限性。相较而言,机器学习算法,如支持向量机、随机森林等,不要求序列的平稳性,可以接受序列本身的统计特征,并将外部环境特征引入模型;然而,处理外部环境对序列的影响需要人工构造特征,进行特征筛选,效率相对较低。

近年来,模型的选用主要集中在神经网络及其优化。相比起传统的机器学习模型,不需要人工构造合适的输入特征,在隐藏层中实现了有效量化因子的寻找,模型预测效果通常优于传统机器学习模型的预测效果。

1.4 研究内容

本项目通过多种方法和多种系统研究,完善已有的高频交易技术数学模型、实战策略等。根据多个品种的 orderbook、trades 等数据,对某品种未来微观价格的变化

做预测，设计相关因子，并进行检验，最后通过线性回归、深度学习等方法进行算法集成，主要应用 LSTM 和 ARIMA 模型进行实盘检验，得出预测结果，并进行对比分析。

1. 处理和清洗原始数据，对数据进行可视化处理。
2. 挖掘高频因子，对短期价格做预测分析。
3. 通过一些算法对高频因子进行集成，加强预测能力
4. 高频交易实盘验证，并开发复盘系统并进行分析。

本研究将从理论和实证两个层面进行探讨。首先，通过文献综述，探讨高频交易技术的概念、特点、优缺点和适用范围。其次，介绍量化对冲基金中常用的高频交易策略和风险管理技术，并探讨其实际应用效果。然后，阐述限价订单簿的概念和应用，并通过案例分析说明其在高频交易中的应用。最后，对常用的量化策略进行设计和实证分析，以验证其有效性和可行性。

1.5 论文创新点

本课题的研究价值在于，通过对高频交易的理论学习和实证分析，实现针对市场价格、成交量、持仓量等真实交易数据的建模分析，应用机器学习实现拟合预测，对模型的实际运用具备实际的指导意义。

1. 综合分析高频交易技术在量化对冲基金中的应用。通过对高频交易技术的深入分析和探讨，综合分析其在量化对冲基金中的应用效果和适用范围，通过多种算法集成使预测结果更为准确。
2. 在理论研究的基础上，设计和实证分析量化策略，进行实盘验证。量化策略是量化对冲基金中的核心竞争力之一，本研究将针对常用的量化策略进行设计和实证分析，将所预测因子在现实中的数字货币市场进行尝试与检验，以验证其有效性和可行性。

第二章 理论模型和研究方法

2.1 高频交易概述

高频交易（High-Frequency Trading, 简称 HFT）是近年来金融领域发展迅速的一种交易方式，以其快速、大量的交易特点广泛应用于量化对冲基金中。高频交易是一种利用先进的数学模型和计算机技术进行高速、自动化交易的方法。这种交易方式可以在毫秒级或微秒级的时间内完成，以追求极短期内的价格差异获取收益。

2.1.1 交易特点

高频交易具有以下几个显著特点：高速交易、大量交易、低延迟、小额交易以及自动化交易。

1. 交易次数多：在相同的周期里，相比中低频交易，高频交易的交易次数最多，也就是交易频率很高。从单个策略来看，交易一个品种每天可以交易几十次到几百次；从单个账户来看，每天的资金使用周转次数从几十次到几百次；如果多策略多品种，每天交易次数多达几千到几万次。

2. 持仓周期短：相比中低频交易，持有资产的时间（也就是周转一次）最短。非要从持有周期上区分高频交易，高频交易持有的部位不会超过一天，也就是日内交易。

3. 盈利稳定：盈利稳定的评价标准就是从每周来看都是盈利的，从每天来看几乎也是盈利的。盈利如此，按日统计，最大的账户资金回撤理论上几乎接近于 0，实际交易中最大回撤应该不超过 1%。

4. 收益率高：如果每天平均盈利千分之一，每年 253 个左右的交易日，年化收益就是 25%。资金量少（千万级）做的好的，年化收益 300% 以上完全不足为奇；资金量大（亿级）做的好的，年化收益超过 100% 的很正常，资金量巨大（十亿级）年化收益 30% 以上算是不错的。

2.1.2 发展背景

近年来，高频交易的迅速发展可以归因于多种因素，其中包括技术创新、市场环境变化、政策法规调整以及对冲基金对高收益的追求。

首先，技术创新是高频交易快速发展的主要原因之一。随着计算机和通讯技术的

迅速发展，交易系统的处理速度和交易数据的传输速度不断提高，从而使得高频交易得以实现。同时，人工智能、机器学习、大数据等技术的发展，也为高频交易提供了更加精准、智能的交易决策和执行手段。

其次，市场环境的变化也推动了高频交易的发展。随着金融市场的全球化和竞争加剧，交易者需要更快、更准确、更高效的交易方式来满足市场需求。高频交易的特点正好符合市场的需求，使其成为了交易者的首选。

此外，政策法规的调整也对高频交易的发展产生了影响。一些国家和地区通过制定相关的政策法规，促进高频交易的发展。例如，美国证券交易委员会曾在 2010 年颁布规定，要求交易平台对高频交易的交易行为进行监管，提高市场的公平性和透明度。

最后，对冲基金对高收益的追求也是高频交易快速发展的一个原因。对冲基金通过使用高频交易来获取更高的投资收益，这也推动了高频交易的快速发展。

总之，高频交易的快速发展离不开多方面因素的推动，其中技术创新、市场环境变化、政策法规调整以及对冲基金对高收益的追求是主要的原因。在未来，随着技术的不断进步和市场需求的变化，高频交易将会不断发展壮大。

2.1.3 策略研究

量化策略是指通过对大量历史数据进行分析，建立数学模型并进行交易决策的一种投资方法。高频交易通过使用量化策略，能够快速捕捉市场机会并迅速执行交易。常见的量化策略包括：市场做市策略、统计套利策略、事件驱动策略、趋势跟踪策略等。

1. 市场做市策略：市场做市策略通过提供买卖价差，实现利润。高频交易在市场做市策略中的应用可以提高市场流动性，降低交易成本。

2. 统计套利策略：统计套利策略是利用历史数据和统计分析方法，寻找市场中价格偏离正常水平的证券并进行交易。高频交易技术可以帮助交易员更快地发现套利机会并实施交易。

3. 事件驱动策略：事件驱动策略关注市场中的重大事件，如公司财报、重大新闻等，并根据事件对市场价格的影响进行交易。高频交易在事件驱动策略中的应用可以在事件发生后的短时间内捕捉市场波动，从而实现交易收益。

4. 趋势跟踪策略：趋势跟踪策略是通过分析市场价格走势，预测未来趋势并根据预测进行交易的策略。高频交易技术可以帮助交易员实时跟踪市场趋势并快速执

行交易。

2.1.4 交易工具

限价订单簿是在现代金融理论中市场微观结构的重要组成部分，常常应用于证券交易领域。限价订单簿是交易所用于显示未执行的限价订单的一种工具。通过限价订单簿，交易者可以看到市场上所有的限价买单和卖单，以及每个价格点上的数量。交易者可以在这些价格点上提交买单和卖单，并等待交易所自动匹配订单。它是订单驱动市场模型的载体，投资者向市场传递不同预期，影响其他投资者调整交易策略，推动证券价格朝新的方向改变。常见订单簿在同一时间会显示买卖各五个价格，分别为：买一、买二、买三、买四、买五；卖一、卖二、卖三、卖四、卖五。其中，未成交的最低卖价是卖一，次低卖价是卖二；而未成交的最高买价是买一，次高买价是买二，以此类推。考虑买入的情况，当存在报价高于卖一时，即时成交；反之，则需等到卖价下跌且更高的买价成交完之后，才有机会成交。

在高频交易中，限价订单簿是一个重要的工具。高频交易者可以利用限价订单簿的信息，快速地识别价格波动和市场趋势，并进行快速的交易决策。通过限价订单簿，高频交易者可以实时了解市场的供需情况和价格波动，从而更好地把握市场机会。

2.2 长短期记忆神经网络 (LSTM)

循环神经网络 (Recurrent Neural Networks, RNN) 是一种非常重要的神经网络模型，用于处理具有序列结构的数据，例如文本、语音、时间序列等等。但是，传统的 RNN 存在着“梯度消失” (vanishing gradient) 和“梯度爆炸” (exploding gradient) 的问题，导致长期依赖 (long-term dependencies) 难以被学习。为了解决这个问题，LSTM (Long Short-Term Memory) 神经网络被提出，并取得了显著的效果。

LSTM 是一种常用的循环神经网络模型，用于处理时序数据，如语音、文本和视频。相比于传统的循环神经网络模型，LSTM 能够更好地解决长期依赖问题，使得网络能够更好地捕捉时序数据中的重要信息。

下面将从经典的 RNN 出发，详细介绍 LSTM 的结构和工作原理，以及如何在模型中使用 LSTM 模块。

2.2.1 循环神经网络 (RNN)

循环神经网络 (RNN) 是一种经典的神经网络模型，主要用于处理具有序列结构的数据。它的基本结构如下图所示：

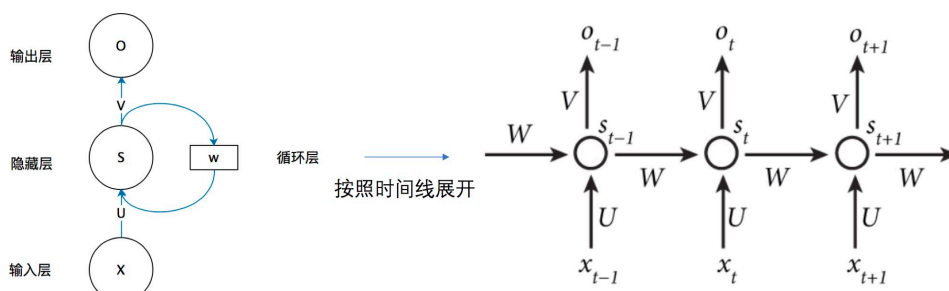


图 2-1 RNN 基本结构

从上图中可以看出，RNN 的结构是由一个循环体 (Recurrent Cell) 和一个输入层 (Input Layer) 组成。输入层接受输入 x_t ，循环体接受上一个时刻的输出 h_{t-1} 和当前时刻的输入 x_t ，经过一定的计算，输出当前时刻的状态 h_t ，并将 h_t 作为下一个时刻的输入。具体地，RNN 的计算公式为：

$$h_t = \sigma(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

其中， W_{xh} 是输入层到循环体的权重矩阵， W_{hh} 是循环体内部的权重矩阵， b_h 是偏置向量， σ 是非线性激活函数，常用的有 \tanh 和 ReLU 等。

虽然 RNN 有着很好的处理序列数据的能力，但是它存在着“梯度消失”和“梯度爆炸”的问题。当序列很长时，反向传播的梯度会随着时间的增加而指数级地衰减或增长，导致长期依赖的信息无法被有效地传递。因此，需要更加复杂的模型来解决这个问题。

2.2.2 LSTM 基本结构

LSTM 是由 Hochreiter 和 Schmidhuber 在 1997 年提出的一种循环神经网络模型，用于解决 RNN 中存在的“梯度消失”和“梯度爆炸”的问题。

LSTM 模型的核心是细胞状态，它是一个长期的内部状态变量，能够记忆历史信息。LSTM 的基本结构如图所示：

在 LSTM 中，每个单元都由四个部分组成：输入门 (input gate)、遗忘门 (forget gate)、输出门 (output gate) 和细胞状态 (cell state)。三个关键的门控结构 (输入门、

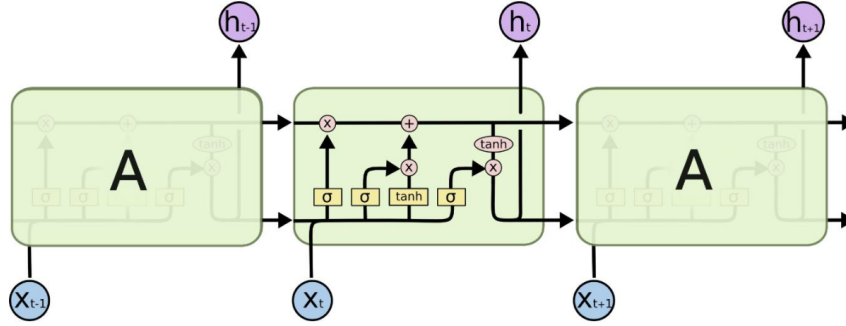


图 2-2 LSTM 基本结构

遗忘门和输出门) 分别都有对应的权重向量, 用于控制信息流的开关。输入门的作用是控制网络需要从当前输入中保留多少信息。

遗忘门的作用是控制网络需要从历史信息中保留多少信息。输出门的作用是控制网络输出的信息量。每个门都有对应的权重向量 (input gate 的权重向量为 w_{xi} , 遗忘门的权重向量为 w_{xf} , 输出门的权重向量为 w_{xo}) 和偏置量 (input gate 的偏置量为 b_i , 遗忘门的偏置量为 b_f , 输出门的偏置量为 b_o)。

2.2.3 LSTM 工作原理

LSTM 模型可以处理时序数据的关键在于它的细胞状态。细胞状态会在每个时间步长中更新, 以捕捉输入序列中的长期依赖关系。

具体来说, LSTM 的工作流程如下:

1. 遗忘门 (Forget Gate)

首先, 需要决定哪些信息需要保留, 哪些信息需要遗忘。遗忘门通过 sigmoid 函数来确定哪些信息需要被遗忘, 计算公式为:

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

其中, W_{xf} 和 W_{hf} 是遗忘门的权重矩阵, b_f 是偏置向量, σ 是 sigmoid 函数。 f_t 表示需要被遗忘的信息。

2. 输入门 (Input Gate)

然后, 需要决定哪些新信息需要被添加到记忆细胞中。输入门通过 sigmoid 函数来确定哪些信息需要被添加, 通过 tanh 函数来生成新的候选信息, 计算公式为:

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

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

其中, W_{xi} 、 W_{hi} 、 W_{xc} 、 W_{hc} 是输入门的权重矩阵和记忆细胞的权重矩阵, b_i 和 b_c 是偏置向量, σ 是 sigmoid 函数, \tilde{C}_t 表示新的候选信息。

3. 记忆细胞更新

接下来, 需要将旧的记忆细胞 C_{t-1} 和新的候选信息 \tilde{C}_t 结合起来更新记忆细胞 C_t 。具体地, 通过将遗忘门 f_t 和输入门 i_t 应用于旧的记忆细胞 C_{t-1} 和新的候选信息 \tilde{C}_t , 得到新的记忆细胞 C_t , 计算公式为:

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

其中, \odot 表示按元素相乘的操作。

4. 输出门 (Output Gate)

最后, 需要决定哪些信息需要输出。输出门通过 sigmoid 函数来确定哪些信息需要输出, 通过 tanh 函数来将记忆细胞转换成输出, 计算公式为:

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

$$h_t = o_t \odot \tanh(C_t)$$

其中, W_{xo} 和 W_{ho} 是输出门的权重矩阵, b_o 是偏置向量, σ 是 sigmoid 函数, o_t 表示需要输出的信息, h_t 表示当前时刻的输出。

综上所述, LSTM 通过引入门机制和记忆细胞, 有效地解决了传统 RNN 中存在的“梯度消失”和“梯度爆炸”的问题, 使得模型可以更好地处理长期依赖的信息。LSTM 通过引入门机制和记忆细胞, 有效地解决了传统 RNN 中存在的“梯度消失”和“梯度爆炸”的问题, 使得模型可以更好地处理长期依赖的信息。在实际应用中, LSTM 被广泛应用于自然语言处理、语音识别、图像描述等领域, 取得了显著的效果。

2.3 ARIMA 模型

ARIMA 模型是一种基于时间序列分析的预测模型, 它通过对时间序列数据进行分析和建模, 预测未来的值。ARIMA 模型可以看作是 AR (自回归) 模型和 MA (移动平均) 模型的结合体, 它将这两种模型中的优点结合在一起, 同时消除它们各自的缺点。ARIMA 模型可以用于处理非平稳时间序列, 其中包括趋势、季节性和周期性等多种非平稳特征。ARIMA 模型可以通过对时间序列数据进行差分操作, 将非平稳序列转化为平稳序列, 然后建立模型进行预测。

2.3.1 自回归模型 AR

AR 部分是指自回归模型，其基本思想是将当前时刻的值与前 p 个时刻的值相关联，其中 p 表示自回归阶数。自回归模型可以用以下公式表示：

$$y_t = c + \sum_{i=1}^p \phi_i y_{t-i} + \epsilon_t$$

其中， y_t 表示当前时刻的值， ϕ_i 表示自回归系数， c 表示常数项， ϵ_t 表示误差项。自回归模型通常适用于具有趋势的时间序列数据，如股票市场的时间序列数据。

2.3.2 移动平均模型 MA

MA 部分是指移动平均模型，其基本思想是将当前时刻的误差与前 q 个时刻的误差相关联，其中 q 表示移动平均阶数。移动平均模型可以用以下公式表示：

$$y_t = c + \epsilon_t + \sum_{i=1}^q \theta_i \epsilon_{t-i}$$

其中， ϵ_t 表示当前时刻的误差， θ_i 表示移动平均系数， c 表示常数项。移动平均模型通常适用于随机波动较大的时间序列数据，如气象预测的时间序列数据。

2.3.3 差分处理 I

I 部分是指差分操作，其基本思想是将非平稳时间序列转化为平稳时间序列。平稳时间序列的特点是均值和方差不随时间发生变化，因此可以更容易地建立预测模型。差分操作的次数表示差分阶数，通常用 d 表示。差分操作可以用以下公式表示：

$$y'_t = y_t - y_{t-1}$$

其中， y'_t 表示差分后的值， y_t 表示原始的时间序列数据。通过差分操作，可以去除时间序列数据中的趋势和季节性变化，使数据更加平稳，有利于建立 ARIMA 模型。

2.3.4 ARIMA 模型的建立

ARIMA 模型的全称是 Autoregressive Integrated Moving Average Model，其数学表达式如下：

$$ARIMA(p, d, q)$$

其中， p 表示自回归阶数， d 表示差分次数， q 表示移动平均阶数。自回归阶数

p 表示当前时刻的值与之前 p 个时刻的值相关联，移动平均阶数 q 表示当前时刻的值与之前 q 个时刻的误差相关联。差分次数 d 用于将非平稳时间序列转化为平稳时间序列。

ARIMA 模型的建立过程一般包括以下几个步骤：

1. 数据预处理：对时间序列数据进行预处理，如去除异常值、缺失值等。
2. 平稳性检验：使用 ADF 检验、KPSS 检验等方法检验时间序列数据是否为平稳时间序列，若不是，则进行差分操作。
3. 自相关和偏自相关函数图（ACF 和 PACF）的分析：根据自相关和偏自相关函数图来确定自回归阶数 p 和移动平均阶数 q 。
4. 模型估计：根据上述分析结果，建立 ARIMA 模型并进行估计，可使用最大似然估计或贝叶斯估计等方法。
5. 模型诊断：对建立的 ARIMA 模型进行检验，包括残差检验、单位根检验等。
6. 模型预测：根据建立的 ARIMA 模型进行预测，并计算预测误差和置信区间等。

ARIMA 模型的预测可以使用以下公式：

$$\hat{y}_{t+h|t} = y_t + \sum_{i=1}^h e_{t+i}$$

其中， $\hat{y}_{t+h|t}$ 表示在 t 时刻预测 $t+h$ 时刻的值， e_{t+i} 表示预测误差。ARIMA 模型的预测误差可以通过计算残差来评估，残差定义为实际值与预测值的差值。

2.3.5 模型优缺点

ARIMA 模型的优点：

1. 适用范围广，可以处理多种类型的时间序列数据。
2. 能够较为准确地预测未来的值，可以帮助用户进行规划和决策。
3. 具有较好的可解释性，可以帮助用户更好地理解时间序列数据。

ARIMA 模型的缺点：

1. 对数据平稳性的要求较高，如果数据不平稳，需要进行差分操作，增加了模型的复杂度。
2. 参数估计需要较多的计算时间，对于大规模的数据集来说，运行速度较慢。
3. 对于异常值和缺失值较为敏感，需要进行数据预处理。

总体来说，ARIMA 模型是一种常用的时间序列预测模型，它可以处理多种类型的时间序列数据，具有较好的预测能力和可解释性。ARIMA 模型的实现需要考虑数

据的平稳性、自回归阶数和移动平均阶数的选择、模型估计和诊断等方面，可以使用统计分析软件或编程语言进行实现。在实际应用中，需要根据数据的特征、领域需求和预测目的来选择适合的模型和方法。

第三章 量化策略与实证分析

3.1 数据集介绍

该项目使用的数据集为校企大创企业方冰宽量化提供的数字货币交易数据集,其包含了各种数字货币交易的数据特征,例如成交量、成交金额、限价订单簿等。这些特征提供了对数字货币交易的全面了解,并为我们提供了探索数字货币市场的机会。数字货币具有很低的基本交易额、方便的交易时间、不易爆仓等特点,因此该数据集为构建数字货币交易的预测模型提供了一个理想的数据源。通过分析这个数据集,可以了解到数字货币市场的一些规律和趋势。这对个人和企业做出投资决策非常有帮助。同时,该数据集可供学习和测试,因此其质量和规模已经得到了保证。

实际运用中,我们首先将限价订单簿数据和交易数据按时间戳排序并将时间戳设为索引。接下来,对于每个限价订单簿数据中的时间戳,我们利用 *pandas* 库中的 *searchsorted* 函数在交易数据中查找该时间戳所在的位置索引。然后,我们可以通过交易数据中该位置索引的前一行数据,确定该时间戳时刻的最新订单簿状态和成交情况。利用该函数对齐数据的好处在于,可以有效减少数据处理时间和空间复杂度,并且避免了循环遍历的过程。因此,该方法适用于大规模数据集的处理。

3.2 模型假设

- 市场是随机的: 高频交易者假设市场价格的变化是随机的,没有规律可循。他们相信,任何价格波动都是由无数个微小的因素组成,这些因素相互作用,并导致价格变化。
- 瞬时的价格变化是无序的: 高频交易者假设瞬时的价格变化是无序的,无法预测。他们相信,短时间内的价格波动是由市场中的交易者的买卖行为引起的,这些交易者的行为是无序的。
- 市场中的信息是不完全的: 高频交易者假设市场中的信息是不完全的。他们相信,无法获取所有的市场信息,因为有些信息是非常难以获取或无法获取的。高频交易者使用可用的信息和历史数据来推断价格的变化趋势。
- 市场反应是瞬时的: 高频交易者假设市场对信息的反应是瞬时的,这意味着价格会立即反映市场的观点和情绪。高频交易者试图从这些瞬时的价格变化中获利。

- 市场是有效的：高频交易者假设市场是有效的，这意味着价格反映了所有可用的信息，并且价格调整反应了市场参与者的需求和供给。高频交易者利用这个假设，使用历史数据来推测未来的价格走势。
- 趋势会持续：高频交易者假设价格趋势会持续一段时间。他们相信，价格趋势是由市场中的买卖行为和市场参与者的心理影响所驱动的。因此，高频交易者会尝试识别价格趋势并跟随其进行交易。

3.3 模型构建

3.3.1 LSTM 模型

本研究基于高频数据集，提取 $voi, oir, cha, voi_2, oir_2$ 作为因子，构建循环神经网络模型 LSTM 实现两方面的预测：1) 预测具体数值；2) 预测涨跌符号。

具体因子计算方法如下：

1) Volume Order Imbalance (VOI)

根据 Shen (2015) 构建的交易量订单流不平衡 (Volume Order Imbalance) 指标 (下称 VOI) 的方法，订单流不平衡测量了在研究的特定时间内，最优买卖价格上委托量的增量之差，反映了投资行为在最优买卖价格上的供需情况，其定义为：

$$OI_t = \delta V_t^B - \delta V_t^A$$

$$V_t^B = \begin{cases} 0 & P_t^B < P_{t-1}^B \\ V_t^B - V_{t-1}^B & P_t^B = P_{t-1}^B \\ V_t^B & P_t^B > P_{t-1}^B \end{cases}$$

$$V_t^A = \begin{cases} V_t^A & P_t^A < P_{t-1}^A \\ V_t^A - V_{t-1}^A & P_t^A = P_{t-1}^A \\ 0 & P_t^A > P_{t-1}^A \end{cases}$$

其中 V_t^B 和 V_t^A 分别是在 t 时刻的买入和卖出交易量， P_t^B 和 P_t^A 分别是最优买卖报价。 δV_t^B 表示买入订单的委托增量。

假设当前的买入报价低于上一期的买入价格，意味着投资者撤单或者订单在的价格上成交了，因此设定 $\delta V_t^B = 0$ ；

假设当前的买入价格和上一期相同，用委托量的增量作为 δV_t^B ；假设当前的买入报价大于上一期买入报价，意味着投资者愿意在更高的价格上买入，价格存在上行的趋势。

同理, δV_t^A 表示卖出订单的委托量增量。将买入订单委托量和卖出订单委托量增量之差作为订单流不平衡指标。

2) Order Imbalance Ratio (OIR)

OIR 是另外一个衡量订单不平衡性质的变量, 能够研究交易者的行为特征。它的统计性质和 VOI 因子相似。我们充分利用盘口各档数据信息, 以期得到更加准确的订单失衡比率 OIR。

VOI 仅衡量不平衡的程度, 仅仅只是买卖委托量的差值, 未考虑到买卖委托量本身的规模大小, 因此其不足以描述市场中交易者的行为。OIR 补充了 VOI 因子, 订单不平衡率帮助我们区分了 VOI 差异大但比率小的情况。例如, 如果当前出价更改量为 300, 而当前要价更改量为 200, 则 VOI 为 100, 这被认为是强烈的购买信号。但此处 OIR 仅为 0.2, 一般大于 0.5 才为较强的买入信号, 因此可以看出原始的买入信号并不是那么强。这里没有考虑买家与卖价订单量之间的比率, 该比率表明了潜在买家和卖家在市场上的相对实力。因此, 我们定义了一个称为订单不平衡率 (OIR) 的新因子。

$$OIR_t = \frac{V_t^B - V_t^A}{V_t^B + V_t^A}$$

其中 V_t^B 和 V_t^A 分别是在 t 时刻的买入和卖出交易量。OIR 为买卖委托量差与其和的比值, 衡量了不均衡程度在其总买卖委托量中的占比。OIR 为正说明市场买压大于卖压, 未来价格趋向上涨, 且 OIR 的比值越大, 其上涨的概率越高, 反之亦然。

3) Bid-Ask Spread

买卖价差作为调节因子, 可以对模型原有因子进行修正。这一因子主要用来衡量流动性, 并调节不同流动性水平下的回归因子, 数学表达式如下所示:

$$S_t = P_t^A - P_t^B$$

。其中 P_t^B 和 P_t^A 分别是最优买卖报价。Bid-Ask Spread 越大, 越难交易。从数学逻辑上看, Bid-Ask Spread 的绝对值越大, 平均价格变化越近于 0, 对交易可能不利。当使用高频数据时, 较大的 VOI 和较小的 $\text{abs}(\text{Bid-Ask Spread})$ 有关。

除了上述三类因子, 本研究还纳入考虑了 Volume Order Imbalance (VOI) 和 Order Imbalance Ratio (OIR) 两类因子的二级结果, 即将公式中相对应的 P 和 V 从订单簿中由买一、卖一的数值替换为买二、卖二的数值。因此, 总共合成了 5 个因子, 作为 LSTM 模型的输入。

建模过程如下所示:

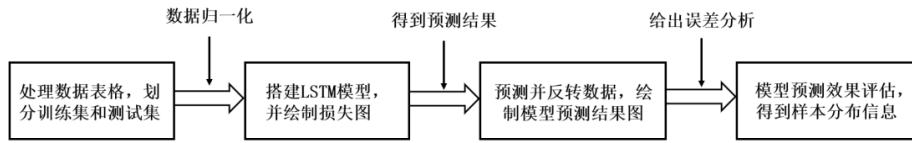


图 3-1 建立 LSTM 模型的流程图

首先对表格数据进行预处理，将每一时刻计算出的各因子与收益单独存储在列表中，并作为模型的输入与输出。而后，为避免使用未来数据进行预测，我们将训练集和测试集所处的时间段分开。具体而言，我们将数组的前 200000 个时间戳作为训练集，并以后 50000 个时间戳作为测试集（使用训练集训练 LSTM 模型，并在测试集上进行预测和评估），同时利用 Z-score 方法对两个数据集各自进行数据归一化处理，即：

$$x' = \frac{x_i - \bar{x}}{\sigma}$$

其中 \bar{x} 为样本平均值， σ 为样本方差。而在搭建 LSTM 基本框架后，对模型超参数进行寻优。综合考虑运算时间和准确率，可以确定各自超参数取值为：

$$hidden_units = [50, 10, 1]$$

$$learning_rate = 10^{-4}$$

$$batch_size = 40$$

$$validation_size = 0.2$$

整个过程中使用 Adam 优化器进行优化，在此过程中记录训练集和测试集的损失并绘制损失图，结果如下所示。

1) 预测具体数值

对于预测具体数值，为了得到最优结果，我们训练了 50 个 epochs，并使用了 early_stopping 方法，patience 设置为 20，即在验证集上表现在 20 个 epochs 内持续下降时，模型停止训练。

由图可知，loss 函数的值随训练次数增多先快速下降，后逐渐下降至大致收敛，说明训练过程较为完整。而由分析可知，用于预测具体数值的模型训练后预测效果较差，呈现预测值绝对值普遍偏小的现象，因此预测效果有待提高。

2) 预测涨跌符号

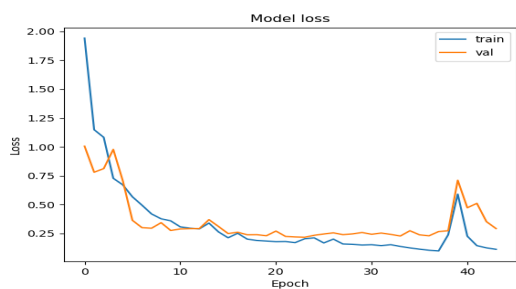


图 3-2 LSTM（预测具体数值）损失函数图像

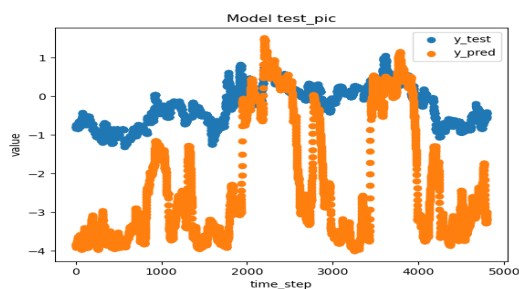


图 3-3 LSTM（预测具体数值）结果

对于预测涨跌符号，我们使用了具有 200 个 epochs 的 LSTM 模型；考虑到该任务下输入特性，为了保证训练的完整性，我们在训练过程中并未使用 early_stopping 方法。

最后得到结果：Test loss = 0.552, Test accuracy = 0.831, 效果相对较好。因此，利用 LSTM 模型预测涨跌符号比预测具体数值的准确率更高，效果更好。

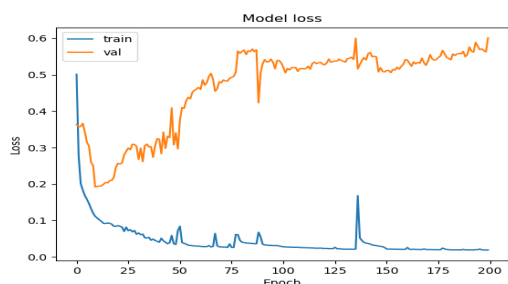


图 3-4 LSTM（预测涨跌符号）损失函数图像

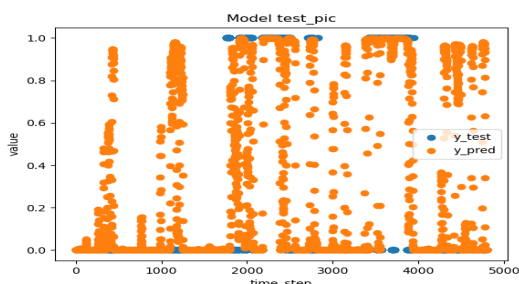


图 3-5 LSTM（预测涨跌符号）结果

3.3.2 ARIMA 模型

本研究将 $t - k_1$ 到 t 时刻的收益作为动量因子、 $t - k_2$ 到 t 时刻的随机波动作为反转因子，用 ARIMA 模型预测 t 到 $t + 2$ 时刻的收益。通过 ACF、PACF 和 BIC 等指标确定动量因子和反转因子的最优时间窗口长短。并用确定的参数构建 ARIMA(2,0,5) 模型预测和进行误差分析。

建模过程如下所示：

1. 模型的识别与定阶

模型的识别与定阶主要是确定 p, q, d 三个参数。首先进行平稳性检验，确定 d 的值。再通过 ACF 和 PACF 函数确定 p, q 的值。



图 3-6 建立 ARIMA 模型的流程图

自相关函数 ACF (autocorrelation function) 主要用于描述时间序列观测值与其过去的观测值之间的线性相关性。 $ACF(k) = \rho_k = \frac{Cov(y_t, y_{t-k})}{Var(y_t)}$ 其中，k 表示滞后期数。

偏自相关函数 PACF (partial autocorrelation function) 描述的是在给定中间观测值的条件下，时间序列观测值预期过去的观测值之间的线性相关性。PACF 与 ACF 的区别在于，ACF 包含了其他变量的影响，而 PACF 则是严格描述了这两个变量间的相关性。我们可以建立 ACF 和 PACF，并分别做出两者的图像。并通过图像按以下方式确定 p, q 参数。其中，截尾指序列从某个时点开始变得非常小。

模型	ACF	PACF
AR(p)	衰减趋于 0	p 阶后截尾
MA(q)	q 阶后截尾	衰减趋于 0
ARMA(p,q)	q 阶后衰减趋于 0	p 阶后衰减趋于 0

图 3-7 确定参数 p, q 的方法

2. 参数估计

通过上述方法定阶具有很强的主观性。本研究在观察 ACF 与 PACF 图像后，得到一定的模型范围区间再通过信息准则函数法，来确定模型的阶数。常用的信息准则函数法有 AIC 准则与 BIC 准则。

最小化信息量准则 (Akaike Information Criterion, 简称 AIC) 是由日本统计学家 Akaike 与 1973 年提出的, 它是拟合精度和参数个数的加权函数: $AIC = 2\ln(L) + 2k$

贝叶斯信息准则 (Bayesian Information Criterion, 简称 BIC) 是 Schwartz 在 1978 年根据 Bayes 理论提出的判别准则, 弥补了 AIC 在样本容量很大时, 选择的模型不收敛于真实模型的不足。 $BIC = 2\ln(L) + \ln(n) \times k$ 其中, L 是在该模型下的最大似然, n 是数据数量, k 是模型的变量个数。通过 BIC 或 AIC 可以以网格搜索的方式绘制热

力图，确定模型最佳的 p, q 组合。

3. 模型检验

本研究主要采取残差检验与回归评价指标检验。残差检验即针对残差进行正态性检验。残差需要满足正态性，服从均值为 0 的正态分布，即白噪声。此时，说明时间序列中有用的部分都已经被提取完毕，只剩下随机扰动。回归指标评价检验中，主要计算预测值于真实值的均方误差、均方根误差、平均绝对误差、 R_2 等定量指标对预测模型进行评估。

3.4 数据分析

本研究探究了时间序列反转效应，对本研究数据建立 ARIMA 模型，以研究过去时刻的收益与随机波动因素对于之后收益的影响。取前 1000 个数据作为研究对象，取在 t 时刻通过对于 $t - k_1$ 到 t 时刻的收益作为动量因子、 $t - k_2$ 到 t 时刻的随机波动作为反转因子，从而用 ARIMA 模型预测 t 到 $t + 2$ 时刻的收益。

首先进行平稳性检验，确认参数 d 的取值 由上图结果我们可以发现，拒绝原假

```
ts = df["benefit"]
ts = ts[0:1000]
ts = ts.dropna()
ts = ts.droptail()
print("差分序列的ADF检验结果为", adfuller(ts))
print("差分序列的ADF检验结果为", adfuller(ts, regression="ct"))
print("差分序列的ADF检验结果为", adfuller(ts, regression="cst"))
print("差分序列的ADF检验结果为", adfuller(ts, regression="nc"))
```

✓ 0.25

差分序列的ADF检验结果为: (-6.4080953583804, 6.441389012988835e-09, 22, 976, {'1X': -3.4370678895881804, '5X': -2.864585868887264, '10X': -2.568349178354273, '18420.843777441896})

差分序列的ADF检验结果为: (-6.73667240122228, 5.520550580567667e-08, 22, 976, {'1X': -3.563075784782453, '5X': -3.414997895317281, '10X': -3.12970132489524, '18420.559291918835})

差分序列的ADF检验结果为: (-6.745051983386933, 3.080536208883974e-07, 22, 976, {'1X': -4.383041117287513, '5X': -3.8384541611720464, '10X': -3.557015215039982, '18418.73090150133})

差分序列的ADF检验结果为: (-6.539267270403979, 6.383157244879406e-10, 22, 976, {'1X': -2.568834586258734, '5X': -1.9412787038611317, '10X': -1.616550690687457, '18421.898957276953})

图 3-8 平稳性检验

设，即不存在单位根，时间序列平稳。

再用 ACF 与 PACF 图确认参数 p, q 的值

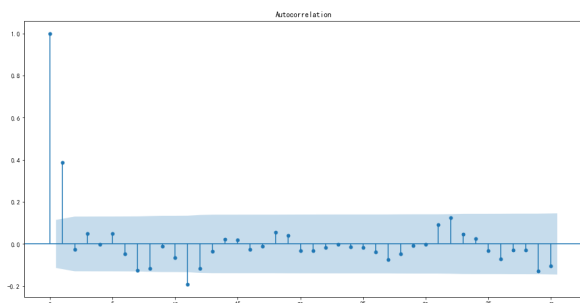


图 3-9 ACF

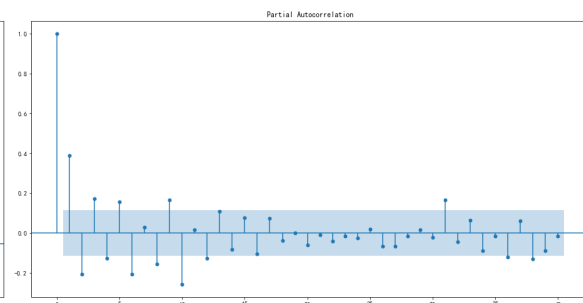


图 3-10 PACF

由上图可发现 ACF 在 1, 11 位置处显著, PACF 在 1, 2, 3, 4, 5, 6, 8, 9, 10, 12, 31, 36, 38 位置处显著。因此选定 p 和 q 的范围分别为 1-38, 1-11, 遍历 38*11 种模型, 选择 BIC 最小者作为结果。

最后构建 ARIMA (1, 0, 1) 模型, 结果如下图所示。

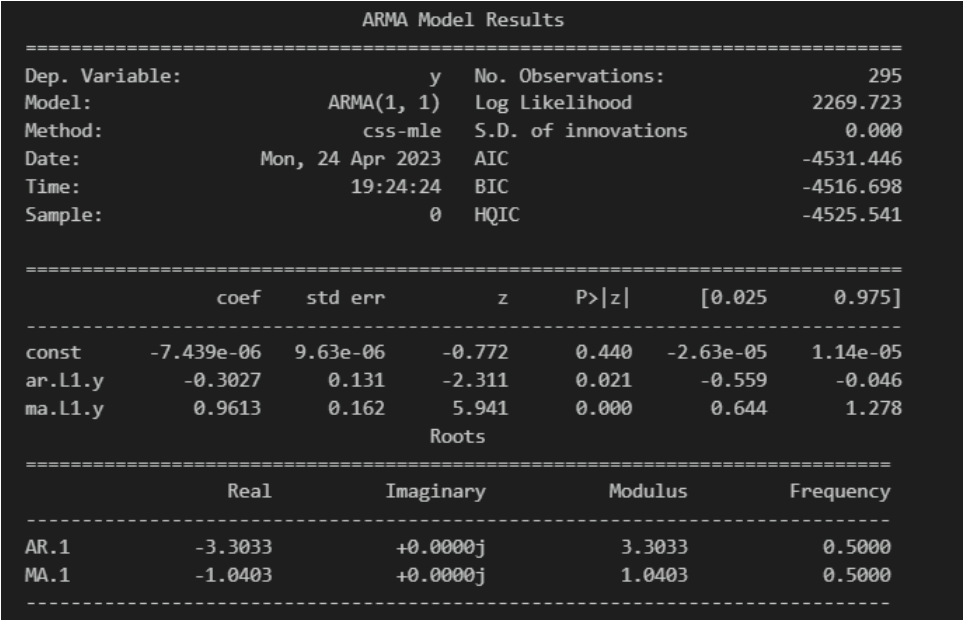


图 3-11 ARIMA (1, 0, 1)

现做出预测值与真实值的对比图像、残差图 (图 12b) 以及残差概率密度分布 (图 12c)。由上图结果可以得出如下结论

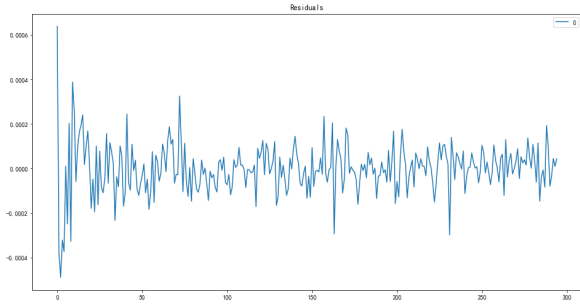


图 3-12 残差

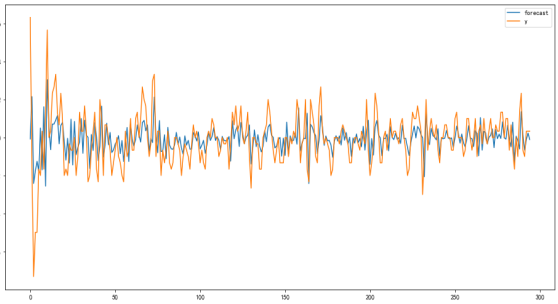


图 3-13 预测

- 预测值与真实值的折线图具有一定的重合度, 计算两者的相关系数为 0.81506057, 为强正相关, 模型趋势拟合较好;
- 除了初始的一段时间, 之后的残差均匀分布在条形宽带中 (由于模型基于过去数据预测, 对于初始数据的预测缺少相应的动量因子和反转因子);

- 残差的概率密度分布近似为均值为 0 的正态分布。利用此模型进行预测，观察该模型十分合理。

最后计算预测值与真实值的均方误差、均方根误差、平均绝对误差、 R_2 等定量指标对预测模型进行评估。结果为均方误差: 0.00000001, 均方根误差: 0.00011396, 平均绝对误差: 0.00007994, R_{square} : 0.30498027 发现预测结果较为良好。

第四章 研究结论与展望

4.1 研究结论

4.1.1 LSTM 模型

考虑到股票数据的序列特性，我们采用了深度学习中的时序模型 LSTM 进行处理并预测相应结果。为避免使用未得到的数据进行预测，我们分离了训练集和测试集所在的时间区间，并结合现实场景的需要，将整体预测任务分为两部分，即股价具体值预测部分和股价涨跌符号预测部分。

在通过超参数寻优和网络训练后，结果显示股价具体值预测任务表现一般，部分结果与真实值偏差较大，此可能与股票市场数据噪声较大，从而干扰预测有一定关系；而在股价涨跌符号预测任务上结果表现较好，测试集上预测准确度较高，此一方面可能与分类任务本身对噪声敏感度较低有关，另一方可能在于在股票市场中，趋势预测任务比值预测任务更加容易有关。

4.1.2 ARIMA 模型

本研究使用了一种新颖的策略，将 $t - k_1$ 到 t 时刻的收益作为动量因子、 $t - k_2$ 到 t 时刻的随机波动作为反转因子，以此来预测 t 到 $t + 2$ 时刻的收益。这样的组合是通过对多项因素进行考虑得出的，通过此方法能够更加全面地考究股票市场的复杂性。

使用 ARIMA 模型进行预测有着广泛的应用，因为它能够从时间序列中提取趋势和周期性。本次研究使用了 ARIMA (2, 0, 5) 模型，其中 2 表示时间延迟为 2 次的自回归项，0 表示没有差分，5 表示时间延迟为 5 次的移动平均项。

通过定阶、参数估计和模型检验三步，本研究最终权衡了精度和简便性，并且得到了收益与动量因子以及反转因子的强关系。研究者还计算了预测值与真实值的均方误差、均方根误差、平均绝对误差、 R^2 等定量指标，证明该模型的预测结果十分可靠。

此外，在本研究中，机器学习算法为股票行业提供了一个全新的视角，并将其与传统的数据分析和投资方法相结合。我们相信该方法可以帮助投资者更全面地认识股票市场的复杂性，更有效地做出投资决策。

4.1.3 模型对比

在本项目中，我们分别采用 LSTM 模型和 ARIMA 模型对股票数据进行预测，得到了相应的结果并由此探索出了模型的不同之处：

- 数据类型 LSTM 模型和 ARIMA 模型都适用于时间序列数据并可以捕捉到序列数据中的依赖关系，但 ARIMA 模型适用于稳定的时间序列数据，比 LSTM 模型的前提要求更高。
- 模型复杂度 LSTM 模型相对 ARIMA 模型更加复杂，需要更多的训练数据和计算资源。LSTM 模型需要大量的参数来学习数据中的复杂模式，可能会导致过拟合的问题。ARIMA 模型则比较简单，只需要较少的参数来建立模型，但可能会出现欠拟合的问题。
- 模型解释性 ARIMA 模型的解释性较好，可以解释模型中每个参数的意义和作用。而 LSTM 模型的解释性相对较差，难以理解模型内部的计算机制和模式。

综上所述，LSTM 模型和 ARIMA 模型各有其优势和局限性，选择何种模型取决于预测任务的性质和要求。在实际应用中，可以根据数据特点和预测目标灵活选择适合的模型。

4.2 不足与展望

尽管 LSTM 和 ARIMA 模型在预测高频交易涨跌方面取得了一定成果，但仍然存在以下不足：

- 数据的实时性：LSTM 和 ARIMA 模型需要大量历史数据进行训练和预测，但在高频交易中，市场信息变化非常快，需要及时响应。如果数据不及时更新，则预测结果可能失去实用价值。
- 算法运行速度：LSTM 和 ARIMA 模型需要进行复杂的计算，运行速度较慢。在高频交易中，时间非常宝贵，每秒钟甚至每毫秒钟都可能产生巨大的利润或损失。因此，算法的运行速度也是非常重要的。
- 模型过度拟合：由于高频交易数据的噪声和波动性较大，LSTM 和 ARIMA 模型可能会出现过度拟合的问题，即模型过于适应历史数据而忽略了新的市场情况。这可能导致预测结果不准确。
- 系统性风险：高频交易在追求短期利润的同时，也面临着系统性风险。例如，市场波动、政策变化、黑客攻击等因素都可能对高频交易产生重大影响，甚至导

致严重损失。

- 难以识别市场变化：高频交易模型可能无法很好地识别市场变化的关键特征，尤其是在非常规的市场情况下。例如，市场恐慌、突发事件、机构交易等因素可能导致市场行为发生变化，高频交易者需要能够快速调整策略来应对。
- 难以预测大事件：高频交易模型可能难以预测大事件的发生和影响。例如，金融危机、地缘政治冲突、自然灾害等因素都可能对市场产生长期和不可预测的影响，高频交易者需要制定更为全面和复杂的风险管理策略。

未来，随着计算机技术和数据科学技术的不断发展，我们可以期待更加先进的高频交易预测算法的出现。例如，使用深度学习和强化学习技术结合市场预测的专业知识来提高预测精度和运行速度。另外，将高频交易数据与其他数据源（如社交媒体数据、新闻事件等）结合使用，可以更好地预测市场走势。因此，未来的高频交易预测算法可能是综合多种数据和算法的复合型算法，能够更好地应对市场的变化和波动。

参考文献

- [1] 文馨贤.(2023). 基于深度强化学习的高频量化交易策略研究. 现代电子技术 (02),125-131. doi:10.16652/j.issn.1004-373x.2023.02.024.
- [2] 周章元. (2022). 基于多因子的市场行业量化投资策略研究 (硕士学位论文, 浙江科技学院).
- [3] 饶瑞, 潘志松, 黎维, 刘松仪, 张磊, 李云波.(2022). 基于深度强化学习的高频交易优化算法. 南京理工大学学报 (03),304-312. doi:10.14177/j.cnki.32-1397n.2022.46.03.008.
- [4] 孙瑞奇.(2016). 基于 LSTM 神经网络的美股股指价格趋势预测模型的研究 (硕士学位论文, 首都经济贸易大学).
- [5] 王苏生, 王俊博, 李光路.(2018). 基于 ARMA 模型的沪深 300 股指期货高频数据收益率研究与预测. 华北电力大学学报 (社会科学版)(03),71-79. doi:10.14092/j.cnki.cn11-3956/c.2018.03.010.
- [6] 黄卿, 谢合亮. 机器学习方法在股指期货预测中的应用研究——基于 BP 神经网络、SVM 和 XGBoost 的比较分析 [J]. 数学的实践与认识,2018,48(08):297-307.
- [7] Aldridge, I. (2013). High-frequency trading: a practical guide to algorithmic strategies and trading systems (Vol. 604). John Wiley & Sons.
- [8] Yao, S., Luo, L., Peng, H. (2018, August). High-frequency stock trend forecast using LSTM model. In 2018 13th International Conference on Computer Science & Education (ICCSE) (pp. 1-4). IEEE.
- [9] Borovkova, S., Tsiamas, I. (2019). An ensemble of LSTM neural networks for high-frequency stock market classification. Journal of Forecasting, 38(6), 600-619.
- [10] Li, Z., Han, J., Song, Y. (2020). On the forecasting of high-frequency financial time series based on ARIMA model improved by deep learning. Journal of Forecasting, 39(7), 1081-1097.

-
- [11] Gencay, R. (1996). Non-linear prediction of security returns with moving average rules. *Journal of Forecasting*, 15(3), 165-174.
- [12] Karazmodeh, M. , Nasiri, S. , Hashemi, S. M. . (2013). Stock price forecasting using support vector machines and improved particle swarm optimization. *Journal of Automation and Control Engineering*, 1(2), 173-176.
- [13] Nelson, D. M. Q., Pereira, A. C. M., de Oliveira, R. A. (2017). Stock market' s price movement prediction with LSTM neural networks. 2017 International Joint Conference on Neural Networks (IJCNN). Presented at the 2017 International Joint Conference on Neural Networks (IJCNN), Anchorage, AK, USA. <https://doi.org/10.1109/ijcnn.2017.7966019>
- [14] S. Selvin, R. Vinayakumar, E. A. Gopalakrishnan, V. K. Menon and K. P. Soman, "Stock price prediction using LSTM, RNN and CNN-sliding window model," 2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI), Udupi, India, 2017, pp. 1643-1647, doi: 10.1109/ICACCI.2017.8126078.

致谢

在项目的研究过程中，我们通过定期的组会分享、理论学习、模型设计和实证检验等一系列努力，实现了对量化交易的方法论研究和实践操作，收获颇丰。在此，我们要向在这过程中给予我们帮助的学校老师和合作企业表示最诚挚的谢意和感谢。

首先，感谢学校提供的校企合作的课题，使我们有机会深入了解量化对冲基金的高频交易技术方法论，开展相关研究。同时，我们还要感谢我们课题项目的合作企业——冰宽量化，在项目前期为我们提供了大量的指导和数据支撑，使我们的研究能够更加深入、精细和实用。

在此，我们还要向各位老师表达最真诚的谢意和感激之情。在研究过程中，老师们不仅为我们提供了宝贵的学术指导和建议，还帮助我们克服了困难和问题，为我们提供了极大的帮助和支持。

最后，我们要感谢本项目组的每一位成员，感谢大家在研究过程中的相互合作和支持，共同推动了本文的顺利完成。相信这段宝贵的经历，将会对我们今后的学习和工作有着长远的积极影响。

再次向各位表达我们最真挚的感谢和敬意！