

分类号 TP311

密级 公开

UDC 004.04

编号

中南财经政法大学

硕士专业学位论文

基于行业的多因子选股和择时策略研究

研究生姓名：张伟钦

校内导师姓名、职称：彭虎锋副教授

申请者类别：非专项计划

校外导师姓名：无

专业学位类别：电子信息

专业名称：电子信息

研究方向：不区分研究方向

入学时间：二〇二〇年九月

二〇二二年六月三日

Research on multi-factor stock selection and timing strategy based on industry

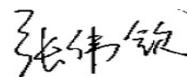
2022.06.03

中南财经政法大学学位论文独创性声明和使用授权声明

学位论文独创性声明

本人所呈交的学位论文，是在导师的指导下，独立进行研究所取得的成果。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的作品。对本文的研究做出重要贡献的个人和集体，均已在文中标明。

本声明的法律后果由本人承担。

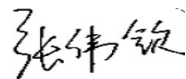
论文作者（签名）：

2022 年 6 月 3 日

学位论文使用授权书

本论文作者完全了解学校关于保存、使用学位论文的管理办法及规定，即学校有权保留并向国家有关部门或机构送交论文的复印件和电子版，允许论文被查阅和借阅。本人授权中南财经政法大学将本学位论文的全部或部分内容编入有关数据库，也可以采用影印、缩印或扫描等复制手段保存或汇编本学位论文。

注：保密学位论文，在解密后适用于本授权书。

论文作者（签名）：

2022 年 6 月 3 日

摘要

随着中国证券市场的快速发展，证券市场日益复杂，传统基于基本面分析的方法在股市中获利越来越困难。伴随着计算机的软硬件快速发展，基于数学和计算机的量化交易逐渐走进了人们的视野。量化选股和量化择时是量化交易两个比较重要的方向，量化选股有很多种实现方式，其中使用较多的是多因子选股，但是近年来的多因子选股论文都把因子与股票未来收益之间的关系看作是相同的，而现实中因子与不同类别股票未来收益的关系不一样。所以基于这点，本文以行业作为股票类别划分的标准，构建了一种基于行业的多因子选股方案。选取的股票能否以正确的时机交易将决定最后股票收益的正负，所以本文结合长期均线、深度强化学习以及行业特点构建了一种新的择时策略，并将其用于分行业选股之后的股票交易阶段，以期尽可能做出正确的交易行为，从而最大化股票的收益。

本文以行业作为股票分类的标准，通过对股市的分析和参考前人研究选取了泛消费行业和生物医药行业作为行业研究对象。并通过综合分析从 A 股中确定了属于泛消费行业和生物医药行业的实验股票。为了尽可能准确刻画股票未来的收益，本文选取了通联量化数据平台上所有免费因子作为初始候选因子，并将因子以行业为单位分两步筛选出构建多因子选股模型的因子。为了分析对比分行业构建多因子选股模型的有效性，本文还加入一个特殊行业：不分行业，即将泛消费行业和生物医药行业的股票合并，并将其以同样的流程构建多因子选股模型。将三个行业使用三种机器学习算法进行训练，并对模型进行评估。通过分析收益发现泛消费行业和生物医药行业大部分时间的收益高于市场平均收益，而不分行业仅有小部分时间的收益高于市场平均收益。所以分行业构建多因子选股模型可以有效的选取未来一段时间内收益较好的股票。

多因子选股模型选取的股票能否以正确的方式交易决定了最终的收益，本文将行业自身的特点、均线指标和深度强化学习的优点结合生成了一种新型择时策略用于交易分行业选取的股票。通过对择时结果的分析，发现泛消费行业采用择时策略的多项收益指标都高于未择时的情况，尤其是大幅度提升了超额年化收益率和夏普比率，显著的提升了相同风险下的收益。而生物医药行业则多个收益指标均有提升。所以本文构建的新型择时策略能够有效做出正确的交易行为，从而最大化收益。

关键词：多因子选股；量化交易；深度强化学习；择时交易

Abstract

With the rapid development of China's securities market, the securities market is becoming more and more complex, and it is more and more difficult for traditional methods based on fundamental analysis to make profits in the stock market. With the rapid development of computer software and hardware, quantitative trading based on mathematics and computers has gradually entered people's field of vision. Quantitative stock selection and quantitative timing are two more important directions of quantitative trading. There are many ways to realize quantitative stock selection. Among them, multi-factor stock selection is used more. However, in recent years, multi-factor stock selection papers have compared factors with The relationship between the future returns of stocks is considered to be the same, but in reality the relationship between factors and future returns of different types of stocks is not the same. Therefore, based on this point, this thesis constructs a multi-factor stock selection scheme based on industry with industry as the standard for stock classification. Whether the selected stocks can be traded at the right time will determine the positive or negative of the final stock returns. Therefore, this thesis combines long-term moving averages, deep reinforcement learning and industry characteristics to construct a new timing strategy and use it for industry-specific stock selection. After the stock trading stage, in order to make the correct trading behavior as much as possible, so as to maximize the return of the stock.

This thesis takes the industry as the stock classification standard, and selects the pan-consumer industry and the biomedical industry as the industry research objects by analyzing the stock market and referring to previous research. And through comprehensive analysis, experimental stocks belonging to the pan-consumer industry and the biomedical industry were identified from the A shares. In order to describe the future returns of stocks as accurately as possible, this thesis selects all free factors on the Tonglian quantitative data platform as the initial candidate factors, and selects the factors to build a multi-factor stock selection model in two steps based on the industry. In order to analyze the effectiveness of building a multi-factor stock selection model for different industries, this thesis also adds a special industry: regardless of industry, the stocks of the pan-consumer industry and the biopharmaceutical industry are merged, and the same process is used to construct a multi-factor stock selection. Model. Three industries were trained using three machine learning algorithms and the models were evaluated. Through the analysis of income, it is found that the income of the pan-consumer industry and the biopharmaceutical industry is higher than the market average income most of the time, but only a small part of the industry's income is higher than the market average income. Therefore, building a multi-factor stock selection model by industry can effectively select stocks with better returns in the future.

Whether the stocks selected by the multi-factor stock selection model can be traded in the correct way determines the final income. This thesis combines the characteristics of the

industry itself, the moving average index and the advantages of deep reinforcement learning to generate a new timing strategy for trading in different industries. selected stocks. Through the analysis of the timing results, it is found that many income indicators of the pan-consumer industry adopting the timing strategy are higher than those of the non-timing situation, especially the excess annualized rate of return and the Sharpe ratio have been greatly improved, and the same has been significantly improved. return at risk. In the biopharmaceutical industry, various revenue indicators have improved. Therefore, the new timing strategy constructed in this thesis can effectively make correct trading behaviors, thereby maximizing returns.

Key Words: multi-factor stock selection; quantitative trading; deep reinforcement learning; timing trading

目录

摘要.....	1
Abstract.....	I
绪论.....	1
第一节 研究背景和意义.....	1
一、研究背景	1
二、研究的意义	2
第二节 文献综述.....	2
一、多因子选股策略的研究现状	2
二、量化择时策略的发展和研究现状	3
三、文献评述	4
第三节 主要的研究内容.....	4
一、 主要的研究内容	4
二、技术路线概述	5
三、内容安排	6
第一章 相关理论知识.....	7
第一节 多因子选股理论.....	7
一、资本资产定价模型	7
二、套利定价理论	7
三、Fama-French 三因子和五因子模型	8
第二节 集成学习和相关算法理论.....	8
一、AdaBoost 算法.....	8
二、XGBoost 算法.....	9
三、LightGBM 算法	9
第三节 强化学习及相关算法理论.....	10
一、强化学习	10
二、深度强化学习 DQN 算法	11
第二章 构建多因子选股模型和择时策略.....	13
第一节 确定研究对象范围和初始股票池.....	14
一、研究行业的确定	14

二、行业代码和所属股票的选取	15
第二节 候选因子选取	15
一、候选因子选取依据	15
二、候选因子分类	16
三、确定候选因子计算周期	17
第三节 数据预处理	18
一、数据处理	18
二、选取构建选股模型的因子	20
三、生成选股模型数据集	22
第四节 构建多因子选股模型	22
一、模型输入输出	22
二、模型训练方式	23
三、模型训练过程	24
四、分类预测结果评估	24
第五节 构建择时策略	26
一、构建长期均线策略	27
二、构建深度强化学习 DQN 择时策略	28
三、构建融合策略	29
四、择时策略评估	29
第三章 选股模型和择时策略实证	32
第一节 实证对象获取和实验环境确定	32
一、数据获取	32
二、确定实验环境	32
第二节 数据预处理	33
一、数据处理	33
二、因子选择	33
第三节 选股模型训练和对比	36
一、LightGBM、XGBoost 和 AdaBoost 模型训练	36
二、训练过程与预测效果对比	37
第四节 构建择时策略与回测验证	43
一、计算长期均线	43

二、训练深度强化学习 DQN 算法	44
三、择时策略实证结果评估	44
第四章 总结与展望	47
第一节 总结	47
第二节 不足与展望	48
参考文献	49
附录 1	52
附录 2	54
附录 3	56
致谢	56

绪论

第一节 研究背景和意义

一、研究背景

随着证券市场的快速发展，证券市场中的股票也越来越多，依靠传统的方法在众多的股票中选出具有投资价值的股票变得越来越困难。同时伴随着科技的发展，人们处理股票的方法也越来越先进，导致预测股票走势的整体难度不断拔高，使得从海量股票中选取收益较好的股票越来越难。面对海量的股票数据，传统的方法已经很难对其进行分析，所以人们开始通过各种机器来协助处理。所以在这种背景下，量化交易渐渐的走入了人们的视野中。

量化交易，通俗的解释就是通过数学模型量化数据，然后使用计算机对量化的数据进行处理，并以计算机发出的交易指令来参与证券市场的交易。它和传统的交易相比，更加的理性，减少了由于人心里因素所带来的干扰。

当前的证券交易市场中，被大多数投资者广泛参与的金融资产就是股票。股票就是企业为了筹集资金，而将自己公司的部分所有权挂在股市上进行“售卖”，交易者通过买入股票从而持有公司一定的股权，并因此成为股东。所有的股东都可以凭借购入的股权去获取股票的股息和红利，所以股票一直以来都被认为具有巨大的盈利潜力。因此如何选取出未来一段时间内具有投资价值的股票就成为广大交易者研究的热点。根据腾讯证券 2020 年中国股民年度行为的研究报告^[20]可以知道，中国的股市中有 75.8% 的交易者持股的时间都不超过 6 个月，即大部分人都是中短线交易者。对于中短线交易者来说，股票的交易几乎可以认定为是一个零和游戏¹。因为在较短的时间内，股票由于自身成长给交易者带来的利润与投入的资金相比，几乎可以忽略不计。所以股票市场中有人盈利就必定有人亏损，而盈利的交易者和亏损交易者的区别就是否正确的判断股票未来一段时间内的走势。所以想要在股市中赚取收益，就必须想办法提高我们判别股票的能力。现实中对股票的判断可以体现在两个方面，一个是选股判断，即选取出未来一段时间内走势优于市场平均走势的股票。另一个是择时判断，即对股票短时间内的价格进行预测，通过预测股票的涨跌来指导交易者的交易行为。

当下对于选股判断，较为成熟的方法是多因子选股。通过研究影响股票走势的因素和股票收益之间的关系，构建出数学模型，然后以当前的数据预测股票未来的收益，从而进行选股。当前多因子选股的大部分研究都是将因子与所有股票未来收益之间的关系看作是相同的，并没有因“股”制宜的研究因子和股票收益之间的关系。所以本文以这个为方向进行了研究，通过行业划分股票类别，构建了一种基于行业的多因子选股模型，以期可以提高多因子选股模型的选股能力。

¹ 零和游戏：指在严格竞争下，一方的收益必然意味着另一方的损失。

对于择时判断,当前的研究既有研究基于传统技术指标来构建择时策略,也有研究基于机器学习相关算法来构建择时策略,但是对于将传统技术指标和机器学习算法进行融合构建择时策略的研究比较少。所以本文结合长期均线和强化学习构建了一种新的择时策略用于股票的择时交易,并根据行业的特点对择时策略进行优化,将其用于分行业构建的多因子选股模型选取的股票。以期可以准确的预测股票的涨跌,从而做出正确的交易行为提高最终的收益。

二、研究的意义

本文主要是对基于行业构建多因子选股模型进行研究,另外通过将传统技术指标和强化学习融合构建新型择时策略,并且构建的策略对股票的历史数据进行模拟回测。所以本文即有理论研究方面的意义,也有实践的意义。

理论意义,首先本文认为因子与股票未来收益之间的关系以股票类别的不同而不同,因“股”制宜的构建多因子选股模型,通过行业将股票进行类别的划分,并以行业为单位选取因子。最后使用选取的因子构建每个行业的多因子选股模型,给多因子选股的研究提供了新的参考。同时本文将传统的技术指标和机器学习进行组合,提取各自的优点建立择时策略,给量化择时研究提供了一种新的思路。

实践意义,本文选取的研究的对象都是我国 A 股中的股票,并且研究数据的年份也是近几年,所构建的选股策略和择时策略均可以直接使用,对广大的交易者具有较强的实践意义。

第二节 文献综述

一、多因子选股策略的研究现状

股票市场是许多国家经济的重要组成部分,对推动世界经济蓬勃发展起着非常重要的作用。股票市场讳莫如深,股票的价格受到各种因素影响,以至于对股票走势的预测十分困难。伴随着数学理论和计算机软硬件的发展,使用数学与计算机结合的方式对股票走势进行预测的研究也逐渐变多,其中探究影响股票走势因素的多因子选股是当下应用最广泛的选股方法之一。

王凯(2017)通过使用逻辑斯蒂算法作为 AdaBoost 算法的基分类器构建了一个新的集成学习算法,并将该算法与多因子选股相结合构建了两种多因子选股模型 M1 和 M2,使用 M1 和 M2 两个模型在 2013 至 2016 年上进行回测²分别取得了 33.36%和 36.7%的年化收益率^[30]。姜加才(2018)将 LightGBM 算法与多因子选股相结合,构建了基于 LightGBM 的多因子选股方案。使用这个方案选出的股票在 2013-2018 年期间取得的总收益为 40.09%,年化复合收益率高达 18.36%的结果^[27]。Zhige Li, Derek

² 回测:根据历史数据来验证交易策略的可行性和有效性的过程。

Yang 等（2019）猜想技术指标与股票未来收益之间的关系会因为股票类别的不同而不同，于是他们从基金的角度出发，以基金选取的股票作为同一个类别的股票进行股票共性的学习，并通过这个共同的属性将 A 股中的股票进行分类^[45]。然后设计一个浅层神经网络用于拟合技术指标与不同类别股票的关系用于优化技术指标，文章中使用这种优化的技术指标在历史数据上回测取得的收益显著高于未使用优化指标取得的收益，以实践的方式证明了这种猜想的正确性。Nguyen Nguyet, Nguyen Dung（2020）等则通过隐马尔可夫模型（HMM）构建一种应用于全球指数的选股策略，文章中挑选出了 5 个对全球指数有影响的因素和 6 个衡量经济的指标，他们认为 5 个因素对于不同的经济指标的影响是不同的，所以他们通过使用 HMM 找出过去时间序列中和当前指标状态相同的时间段，并重新分析出每个因素对经济指标的影响权重，最后根据分析的权重用于选股。文章最后的结果显示在 2005 年 1 月至 2017 年 9 月使用 HMM 选股达到了 210% 的收益率，而使用等权重选股情况的收益率只有 70% 的收益率，所以通过 HMM 有效的提升了选股能力^[39]。Huotari Tommi, Savolainen Jyrki 等（2021）等人则把深度强化学习应用于选股，并最后通过实证的方式证明深度强化学习可以运用于金融领域，但是文中也指出深度学习构建的策略是一个具有高风险的策略^[40]。袁晨光（2021）则通过将朴素贝叶斯、逻辑回归、支持向量机和 LightGBM 等算法构建硬、软集成学习算法用于多因子选股^[5]。文章最终的结果是构建的多种集成学习算法在回测中获取的收益都显著的高于沪深 300 指数在回测中取得的收益。

通过上述量化选股的文献可以知道，量化选股包含许多种方法，但是所有的方法中都围绕一个核心，即选取出未来一段时间内具有较好收益的股票。其中多因子选股是一个比较成熟的选股方法，通过分析大量的实证可以发现多因子选股是一种有效选取出未来收益较好的股票的方法。

二、量化择时策略的发展和研究现状

华宇（2019）将 BP 神经网络运用于构建量化择时策略，将股票的日数据输入到 BP 神经网络中进行训练，并以股票的收盘价作为输出^[24]。最后文中指出基于 BP 神经网络预测股价的准确率比传统技术因子的准确率要高，所以基于神经网络预测股价并用于构建择时策略是一种有效的方法。何路（2020）通过情绪指数来构建择时策略，将情绪指数引入到交易阶段，依据情绪指数取值的不同来预测股票在未来一小段时间内的涨跌，进而产生交易信号^[15]。文章中最后的结果是引入情绪指数显著的提高了模型的择时能力。常太星（2020）通过将隐马尔可夫模型（HMM）来构建择时策略，文中指出 HMM 在市场动态判定方面具有较高的精确性，能够有效对市场反转周期做出判断，进而有效的判断当前股票所处市场的状态，最后依据当前市场所处的状态产生准确的交易的信号^[14]。董焕彬（2020）则将多智能体强化学习引入择时交易^[10]。文章中指出多智能体强化学习借助其强大的数据处理能力，能更好的处理金融市场的

海量数据。在择时方面，能够提升交易信号的信息提取能力，从而取得超额收益。Koshiyama Adriano, Firoozye Nick 等（2021）则将对抗神经网络（GAN）用于择时交易，文章中通过 GAN 对多种交易策略进行校准和组合从而在交易过程中获得优势，进而提高交易的收益^[36]。Budiharto Widodo（2021）等将长短期记忆网络（LSTM）来预测印度尼西亚在新冠疫情期间股票的价格，文章通过实验指出他们构建的 LSTM 方法适用于短期（一年）数据的预测，通过这种方法的预测准确率高达 94.57%^[38]。

通过上述的研究可以发现，通过对历史股票价格、市场动态等因素的分析并以此对未来股票价格和市场动态进行预测，从而做出合理的交易行为的择时策略是可以显著的提升收益。

三、文献评述

通过翻阅国内外文献可以看到量化选股和量化择时策略经过长期的研究和实证已经使得这两个领域积累了大量的研究成果，对这些研究进行深入分析，可以做出如下总结：

首先是量化选股方面的研究，近年来的量化选股研究大多数都是通过分析影响股票走势的因素从而来预测股票未来的走势。其中多因子选股相对于其他的选股策略将更多的因素纳入了研究的范围，从而更加精准的把握了股票未来的走势。但是值得注意的是许多的多因子选股研究中，都是将因子与股票未来收益之间的关系看作是相同的，并没有因“股”而异，因“股”制宜的使用因子构建多因子选股模型，所以本文在上述研究的基础上对因子进行区别的对待，以行业来划分股票类别，构建基于行业的多因子选股模型。

其次是量化择时方面的研究，采用机器学习的方法进行预测股价的预测是当前较多人使用的方法，通过机器学习对于历史股票走势进行拟合从而提升预测股票未来价格的概率并以此决定如何交易从而提升收益。而传统的技术指标用于择时的方法取得较好收益并取得成效也有人研究。但是结合传统指标和强化学习的优势，融合形成新的策略用于择时却鲜有人研究，故本文通过传统均线指标和强化学习 DQN 算法构建全新的择时策略并应用于分行业构建的多因子选股模型选取的股票，以探究这种新型择时策略的有效性。

第三节 主要的研究内容

一、主要的研究内容

基于上述的研究背景以及当前研究的现状，本文确定了选股和择时策略两个方面的研究内容：

- 1.以行业作为股票的分类标准，为不同行业的股票选取各自的因子并通过使用 LightGBM、XGBoost 和 AdaBoost 等三种算法对选取的因子和股票收益之间的关系进

行拟合构建多因子选股模型。

2.使用长期均线和深度强化学习 DQN 算法构建一种新的择时略，同时分行业对择时策略进行优化，并将择时策略对分行业构建的选股模型选取的股票进行择时交易。其中 DQN 算法主要是用于规避短期内股票回调，从而最大化收益。

综述所述本文和其他人的主要区别是因“股”制宜，对不同行业的股票的分别选取各自的有效因子，并以各自的有效因子构建多因子选股模型。同时对多因子选股模型选取股票之后如何进行交易的问题做出了回答，本文结合传统技术方法结合强化学习生成了一个新型择时交易策略。

二、技术路线概述

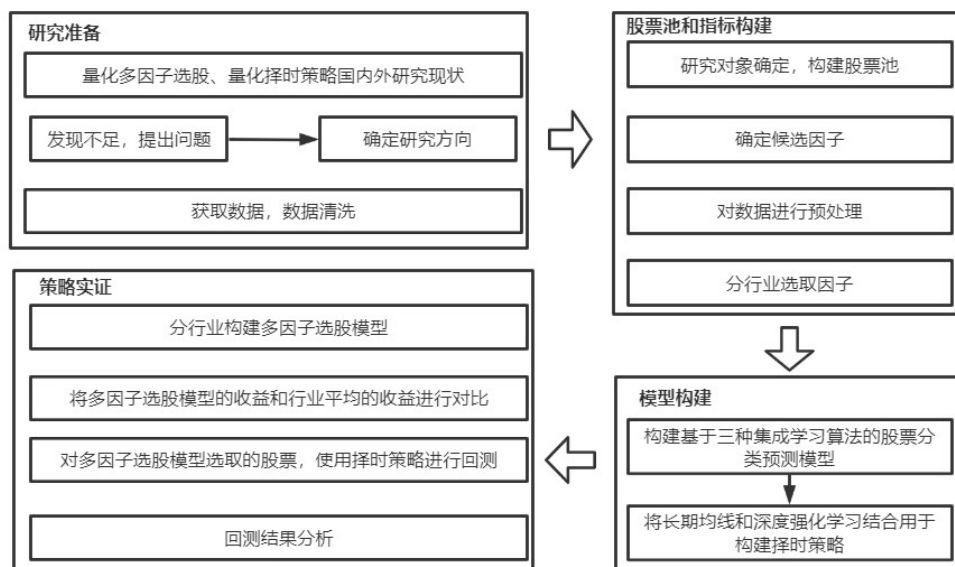


图 0-1 本文技术路线

1.研究准备。首先阅读国内外量化多因子选股和量化择时相关的论文，把握当前研究现状并对当前研究进行分析思考确定出本文的研究方向，然后收集相关实验数据并对数据进行清洗。

2.股票池和指标构建。根据当前股市的实际情况和以及现实的实验条件确定研究的行业对象并以此确定初始股票池，然后基于现实条件选取候选实验因子。最后对确定的股票数据和因子数据进行预处理。

3.模型构建。将实验对象通过三种 LightGBM、XGBoost 和 Adaboost 三种算法分行业分别训练构建多因子选股模型。在多因子选股模型选取股票之后对选取的实验股票通过使用长期均线和深度强化学习进行结合训练生成择时策略。

4.策略实证。对选股模型策略进行技术指标和实际回测的实证，同时也对选取的股票使用构建的择时策略进行回测实证用以验证多因子选股模型是否有效的选取出

收益较好的股票以及择时策略是否能够有效的做出正确的交易行为从而提升最终的收益。

三、内容安排

本文以发现问题、提出解决问题的思路、设计研究的方法、实证和总结的思路来写作，所以论文的主要可以分为如下五个部分：

绪论。这部分内容主要是交代研究的背景，并着重介绍这个背景下，国内外当前研究的现状，并对当前研究内容进行总结，从而引出本文研究的内容。最后就是交代本次论文的写作的大致框架以及论文的组成结构。

第一章，相关理论知识。首先介绍本文选股模型相关的理论，主要是介绍多因子选股相关理论，并以此解释为什么多因子选股有效。其次对着多因子选股使用的集成学习算法进行介绍，最后就是对择时策略中使用的强化学习以及涉及的相关算法进行介绍。

第二章，阐述如何构建多因子选股模型和择时策略。首先确定本次研究的具体行业对象，确定初始股票池。其次确定初始候选因子库，并结合因子选择方法来剔除无效因子，然后分行业使用 LightGBM、XGBoost 和 AdaBoost 等三种集成学习算法构建股票分类预测模型，最后介绍如何构建基于长期均线和深度强化学习 DQN 算法融合的择时策略。

第三章，对多因子选股模型和择时策略的实证。主要是根据第二章中阐述的方法进行实验，主要包括股票的筛选，确定满足实验条件的股票；因子的选择，确定构建多因子选股模型所使用的因子；构建基于行业的多因子选股模型，将多因子选股模型的结果进行对比分析；根据选股模型选取的股票，计算每个股票长期均线的取值，同时也对每个股票使用深度强化学习算法训练得到择时模型，最后将两者进行融合，生成一种全新的择时交易策略。

第四章，总结与展望。这个章节主要是对本文的研究的内容和实验得出的结果进行总结评价以及对未来本研究还可以尝试的方向进行介绍。

第一章 相关理论知识

本章将着重介绍本次研究所涉及到的相关理论，首先介绍的是选股理论，主要是介绍多因子选股相关的理论；然后介绍本文构建多因子选股模型所使用的机器学习算法的基本原理；最后介绍择时策略中使用的深度强化学习和 DQN 算法的基本原理。

第一节 多因子选股理论

影响股票走势的因素多种多样，通过研究不同因素与股票未来收益之间的关系，以数学模型构建因子与收益之间关系的多因子选股模型是当下量化交易领域最成熟的选股模型之一。量化选股模型从最初的投资组合模型、资本资产定价模型（CAPM）和套利定价理论（APT）等理论发展到现在的多因子选股模型，逐渐变得成熟。现实中，由于交易者的学识和经验不同，对市场的理解也千差万别，因此构建的多因子选股模型千奇百态。

一、资本资产定价模型

资本资产定价模型（Capital Asset Pricing Model, CAPM）是多因子选股模型的基础理论之一，它是由夏普（Sharp）等学者通过研究 Markowitz 的均值-方差理论，并在此基础上进行改进，从而发明了可以应用到股市中的投资模型。它解释了均值方差理论中市场平衡状态的来源，并且使用数学模型将投资收益与投资所面临的风险之间关系表达出来，极大的减少了资产组合运算的复杂度。投资回报的数学表达式如公式（1-1）所示：

$$E(r_i) = r_f + \beta_i \cdot (E(r_m) - r_f) \quad (1-1)$$

其中： $E(r_i)$ 代表第 i 个资产的预期回报率； r_f 代表无风险利率； $E(r_m)$ 是市场 m 的预期市场回报率； $(E(r_m) - r_f)$ 为市场 m 的风险溢价； β_i 是资产 i 的系统性 Beta 系数，即风险系数^[17]。

二、套利定价理论

套利定价理论（Arbitrage Pricing Theory, APT）是罗斯（Ross）通过进一步研究 CAMP 模型，并在这个模型的基础上进一步延伸而产生理论。该理论和资本资产定价模型理论有很多相同的基础，比如它们都是市场平衡下的理论。但是套利定价理论相比于资本资产定价理论来说，它的假设条件要少，它将收益率形成过程中的多因子模型作为基础，并且套利定价理论有两个非常重要的论述：

1. 证券收益率与一组代表证券收益率的因子线性相关。

2. 无风险套利的机会诞生于处在不平衡状态中的市场，但是无风险套利机会也会倒逼处在不平衡状态中的市场走向平衡状态。

三、Fama-French 三因子和五因子模型

(一) Fama-French 三因子模型

在夏普的资本资产定价模型中股票的期望收益只与市场的系统风险 B 系数相关，但是在随后的一系列研究中发现有股票的收益还与账面市值比、市盈率倒数等一系列指标有关。Fama 和 French 于 1992 通过研究发现市值因子、账面市值比因子可以解释大部分股票价格的变动，并由此提出了基于市值因子、账面市值比因子和市场风险因子的三因子模型，其公式如 (1-2) 所示：

$$R_{it} - r_{ft} = \alpha_i + \beta_{1i}(R_{mt} - r_{ft}) + \beta_{2i}SMB_t + \beta_{3i}HML_t + \epsilon_{it} \quad (1-2)$$

其中 R_{it} 代表在时间 t 时资产的收益率， r_{ft} 代表无风险收益率； R_{mt} 则为在时间 t 时市场收益率； SMB_t 、 HML_t 则分别代表市值因子和账面市值比因子在时间 t 的收益率， β_1 、 β_2 和 β_3 分别为三个因子的系数， ϵ_{it} 为残差项， α_i 为截距。

(二) Fama-French 五因子模型

如果三因子模型中的三个因子可以完全解释股票的收益，那么公式 (1-2) 中 α_i 的真实值应为 0，但是在三因子模型提出后的 20 多年里许多学者都通过实证分析发现 α_i 的值显著不为 0，所以三因子模型并不能解释所有股票的收益。后来，Fama 和 French 发现除上述因子之外，盈利水平因子以及投资水平因子等因子也能够带来股票的收益，所以 Fama 和 French 于 2015 年提出了五因子模型，其公式如 (1-3) 所示：

$$R_{it} - r_{ft} = \alpha_i + \beta_{1i}(R_{mt} - r_{ft}) + \beta_{2i}SMB_t + \beta_{3i}HML_t + \beta_{4i}RMW_t + \beta_{5i}CMA_t + \epsilon_{it} \quad (1-3)$$

其中 RMW_t 则是盈利因子在时间 t 时的收益率，而 CMA_t 则是投资因子在时间 t 时的收益率， β_4 和 β_5 则是盈利因子和投资因子的系数。

随着 Fama-French 五因子模型的提出，许多学者在此基础上进一步研究，不断地加入新的因子来解释股票的收益，因此多因子模型渐渐成为当前研究的热门。多因子选股模型的主要思路是通过分析和实证找出可以有效解释股票收益的因子，然后利用这些因子来选取股票。

第二节 集成学习和相关算法理论

一、AdaBoost 算法

AdaBoost 全称是自适应增强(Adaptive Boosting)，它是一个非常优秀的集成学习算

法，可以主动适应每个基分类器训练产生的误差，然后将简单的基分类器进行提升误差率，最后使用不同权重将基分类器进行组合得到一个强的分类器。其训练的整个过程为如下所述：

- 1.首先，算法会对初始数据的一部分进行训练，得到第一个训练的好的基分类器。
- 2.其次，对第一次训练出现错误的数据与新的数据结合起来进行二次训练，并得到一个基分类器。
- 3.然后，将第一次和第二次训练中得到的错误数据组合起来进行第三次训练，这样得到第三个基分类器。
- 4.最后，将各个基分类器通过一定的权重组成一个强分类器。

通过对 AdaBoost 的训练描述我们可以知道，算法进行训练的时候会将训练样本分成多个部分，然后对分类错误的数据进行多次训练，从而提升分类器的分类能力，最后对不同的分类器使用不同的权重进行组合得到新的分类器。

二、XGBoost 算法

XGBoost (eXtreme Gradient Boosting) 是一款优秀的算法。在 XGBoost 中为了获取到使损失函数降到最低的基分类器，使用了二阶导数。XGBoost 算法主要是将目标函数进行了改写，通过改成泰勒公式，并且利用泰勒的二阶展开公式来降低训练时候的损失函数，并且以此降低了过拟合的风险。不过事物都是具有两面性的，这个改进有一些缺点，正如上述列举的众多公式一样，它的运算的复杂性非常高，计算量十分巨大，所以运行这个算法耗时相对来说较为漫长。

三、LightGBM 算法

LightGBM (Light Gradient Boosting Machine) 是微软提出的一种新型算法，它从原理层面来说和 XGBoost 类似，其主要的基模型是决策树模型。但将其与 XGBoost 进行比较，可以发现他比 XGBoost 有两个明显的优点：1.LightGBM 的训练速度比 XGBoost 算法的训练速度快；2.LightGBM 算法训练时占用的内存要比 XGBoost 训练时占用的内存少。LightGBM 算法主要是进行了如下几个地方的改进：

- (1) Histogram 算法：它的主要思路是根据直方图来实现最优点的分割^[9]。
- (2) GOSS 算法：它的主要思路是通过合理的去除冗余样本来提高运算速率。主要采用的是一个基于梯度的思想，通过单边采样算法保留在训练过程中产生较大误差的样本，然后集中对这部分误差较大的样本进行训练，使得整体的误差逐渐减少。
- (3) EFB 算法：它的主要思路是通过合理的去除冗余特征来提高运算速度，采用特殊的特征合成算法，将样本中互斥的特征进行融合，从而形成新的特征，以此达到降低整体维度特征的目的，从而有效的提升了训练的速度。

综合上面的解释，我们可以知道，LightGBM 其实就是对运行的内存、运算的数据和训练的特征等三个角度进行优化改进才得以使得算法具有占用内存小和训练

速度快等优点。

第三节 强化学习及相关算法理论

一、强化学习

强化学习属于机器学习范畴，它受到心理学的启发，通过模拟人类学习知识的方式对未知的事物进行学习。它主要由智能体（Agent）、状态（State）、环境（Environment）、动作（Action）、奖励（Reward）组成等五个部分组成。智能体³通过当前的状态做出动作，使得智能体从当前状态转移到另一种状态，并获取到新状态的奖励，智能体通过获取到的奖励来优化做出的动作，使得智能体形成一个达到最大奖励的策略。智能体与环境的交互主要是以状态、动作和奖励来进行，这种交互方式和人类学习时的交互方式有异曲同工之处。

当前强化学习主要分为基于模型（Model-based）和不基于模型（Model-free）两种，两者最大的区别是强化学习是否根据环境的信息重新建模。如果对环境信息重新建模，那就是基于模型的强化学习，如果没有对环境建模那就是不基于模型的强化学习。

（一）基于模型的强化学习（Model-Based）

基于模型的强化学习指的是智能体会通过与环境的交互逐渐学习到真实环境中的规则，比如状态的改变会得到怎么样的奖励，通过这些规则对真实的环境进行建模，一旦完成了模拟环境的构建，智能体与真实环境的交互会急剧减少。基于模型的强化学习需要对环境进行建模，但是大多数的现实场景都非常复杂，难以对其进行建模，比如股票市场，几乎无法对其进行建模，因为影响股票价格的因素错综复杂。

（二）不基于模型的强化学习（Model-Free）

不基于模型的强化学习指的是智能体不需要对环境进行建模，通过与环境的交互中获取到学习的数据，并通过这些数据对做出动作的策略进行优化，从而使得策略可以最大化的获取奖励。不基于模型的强化学习因为不用对环境进行建模，直接从与环境的交互中进行学习，相比于基于模型的强化学习，它的适用范围更加广泛。

本文是将强化学习应用在股票市场，由于股票市场非常的复杂，因此要对其进行建模是一项几乎不可能的事情，所以本文使用的是不基于模型的强化学习，并使用深度强化学习 DQN 算法来构建策略。

³ 智能体：强化学习中行为的决定者。

二、深度强化学习 DQN 算法

深度强化学习 DQN 算法是将深度学习和强化学习 Q-learning 算法进行结合，使用神经网络弥补了 Q 值表爆炸的问题。

（一）Q-Learning 算法

Q-Learning 算法是一种经典的强化学习算法，更准确的说，是一种关于策略的选择方式。通过强化学习描述，可以发现，强化学习的核心就是如何构建一个恰合时宜的策略，使得智能体在环境中的每一次交互都可以获取到尽可能多的奖励。Q-learning 是一种基于值的算法（Value-based）⁴，其主要是根据值的大小来选取动作，并根据奖励来更新策略，具体来讲就是分为如下几步：

首先，算法输入当前的状态 S ，根据计算公式计算每个动作的 Q 值。

其次，选取 Q 值最大的动作作为输出动作。

然后，根据选取的动作与环境进行交互，并获得执行动作后的奖励。

最后，将奖励用于更新计算公式，使策略朝着奖励最大化的方向迭代。

Q 值的更新如公式（1-4）所示：

$$Q'(s, a) = Q(s, a) + \alpha(r + \gamma \max_{a'}(s', a') - Q(s, a)) \quad (1-4)$$

其中 s 为状态 State， a 为采取的行为 Action， α 参数用来表示新的值对更新后值所造成的影响大小， r 在状态 s 下采取动作 a 后获得的奖 Reward， γ 也是一个折扣值，即用来减小新值的影响的值。其中 α 和 γ 的范围都在 0-1 之间。

Q-learning 算法的主要思想是将当前的状态（State）和动作（Action）之间构成一张 Q 值表，所有的状态和动作都存储在这张 Q 值表中。当环境比较简单的时候，状态的量不多，使用 Q 值表进行存储也不会对内存造成很大的负担，但是当状态的数量趋于无穷多的时候，如果还使用 Q 值表来存储状态和动作，将会对计算机内存造成巨大的负担，甚至有时当下最大的内存也无法存储 Q 值表。所以为了解决这种 Q 值爆炸的问题，对 Q-learning 算法做了一些改进，通过使用深度学习和 Q-learning 的进行结合，将 Q 值使用神经网络来预测，从而解决存储 Q 值的问题，这样诞生了强化学习的另一种算法 DQN 算法。

（二）DQN 算法

DQN 是基于 Q-learning 算法，并对不足的地方加以改进而延伸出来的算法。通过神经网络将 Q 值表更新的问题转化成为函数拟合的问题，从而解决 Q 值存储的问题。将深度学习和强化学习进行结合需要解决如下几个问题：

1. 深度学习的训练需要大量的带标签的训练数据，但是强化学习只能提供奖励

⁴ 基于值的算法：即算法为给每个动作生成一个值，智能体根据值来决定下一步的行为。

(Reward) 作为反馈。

2.深度学习的训练数据独立，但是强化学习的前后状态是有关联的，前面的状态有可能会对后面的状态产生影响。

3.使用非线性的网络去拟合值会出现不稳定情况。

为了解决上述问题，DQN 算法采用了如下的几种办法：

1.将奖励和状态动作等构造深度学习需要的训练数据，即将奖励作为标签，状态和动作作为训练特征。

2.构造一个用于存储训练数据的容器，首先将初始数据放入到容器中，等待容器中有足够多数据的时候，随机抽取出一部分数据用于训练，这样就达到了训练数据样本独立的目的。

3.通过构造两个神经网络，一个是训练时的网络输出估计的 Q 值，实时更新参数，一个是作为输出 Q 真实值 (Q-True) 的网络，不实时更新参数。估计的 Q 值不断的向 Q 真实值方向迭代。

由上述描述可以知道 DQN 算法的损失函数为如下公式 (1-5) 所示：

$$L(\theta) = E(TargetQ - Q(s, a; \theta))^2 \quad (1-5)$$

其中 θ 是网络参数，而 Q 和 $TargetQ$ 分别是训练网络和目标网络。

第二章 构建多因子选股模型和择时策略

本章主要是阐述如何构建基于行业的多因子选股模型和择时策略，而实证部分则在第三章。构建基于行业的多因子选股模型和择时策略主要分为如下五个部分：

第一部分，确定研究对象范围以及确定初始股票池。当前的多因子选股研究中都认为因子与所有股票未来收益的关系都相同，但是在 Zhige Li, Derek Yang 等（2019）的论文中发现^[45]，股票类别不同，因子与其未来收益的关系也不同。所以本文对因子以因“股”制宜的方式构建多因子选股模型，即基于股票类别构建多因子选股模型。对于股票类别的划分常用的方法有两种：1.以行业生命周期的不同阶段对股票进行分类，将处于相同行业生命周期阶段的股票看作是相同类别的股票；2.以行业对股票进行分类，对属于同一个行业的股票看作是相同类别的股票；考虑到数据获取的难易程度和实验的复杂度，本文以行业作为股票类别的分类标准，将隶属于同一个行业下的股票看作是相同类别的股票，把基于股票类别构建多因子选股模型的想法通过分行业构建多因子选股模型来实践验证。根据前人的研究以及对股票市场的分析，本文选择了泛消费行业和生物医药行业两个行业作为研究行业，并将两个行业包含的 481 只股票作为初始股票池。

第二部分，候选因子选取。本部分内容主要是确定初始候选因子，为了尽可能准确的刻画股票未来的收益，本文选取了通联量化数据平台上所有可以使用的因子，共计 232 个因子作为初始候选因子。

第三个部分，数据预处理。数据预处理主要包含数据处理和选取构建选股模型的因子以及生成选股模型数据集三个部分，其中数据处理主要是股票的退市和异常检查、股票数据的标签化以及股票和因子数据的标准化处理。选取构建选股模型的因子部分则是对初始候选因子进行初步相关性筛选和最终因子筛选。生成选股模型数据集则是将标准化的因子数据和股票标签值组合形成模型可以使用的数据集。

第四个部分，构建多因子选股模型。这个部分将阐述如何分行业使用 LightGBM、XGBoost 和 AdaBoost 等三种算法构建多因子选股模型。

第五个部分，构建择时策略。这个部分主要是阐述了如何分行业构建基于长期均线和深度强化学习的新型择时策略。

构建基于行业的多因子选股和择时策略建模相关流程如图 2-1 所示：

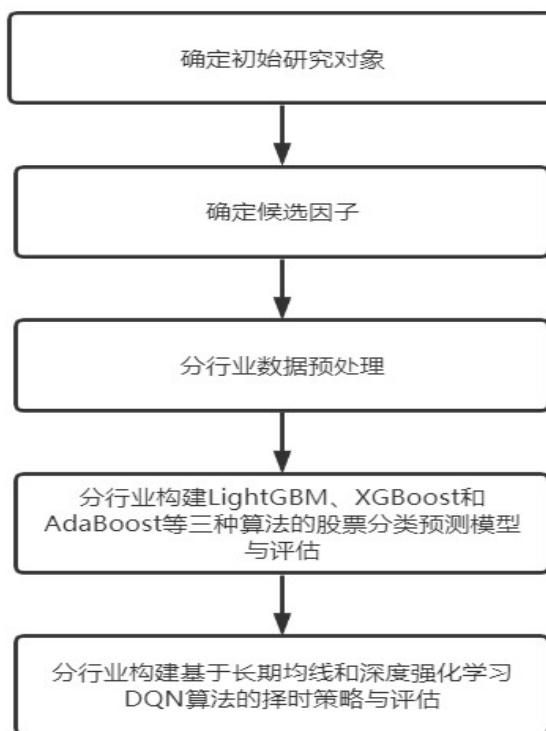


图 2-1 基于行业的多因子选股和择时建模流程

第一节 确定研究对象范围和初始股票池

一、研究行业的确定

在 Zhige Li, Derek Yang 等（2019）的论文^[45]中指出，因子与不同类别股票未来收益的关系并不相同，所以本文从这个观点出发，因“股”制宜以同一个类别的股票为单位构建多因子选股模型进行研究。对于股票类别的确定，当前主流的方法主要有如下两种：

1.以行业生命周期的不同阶段为依据，对股票进行分类。每个行业的生命周期内会经历多个不同的发展阶段，处于不同生命周期阶段的行业各有各的特点，比如刚刚创立的新兴行业，处于行业生命周期的初创期，具有高风险、低收益的特点。而处于成长阶段的行业又具有高风险、高收益的特点。处在不同行业生命周期内的公司或多或少都具有该行业生命周期阶段的特点，而公司的这些特点会直接表现在股票上，所以可以通过行业生命周期的不同阶段来划分股票的类别。

2.以行业为分类依据，对股票进行分类。对于隶属于相同行业的公司来说，毫无疑问都具有该行业的特点，同样公司的这些特点也会表现在股票上，所以以此作为股票分类的依据是一种可行的方法。

两种分类方法各有各的特点，如果从适用性的角度来考虑，那么选取第一种方

法相对更为合适，因为使用第一种方法分类的股票包含了更多深层次的共性，所以以此构建的多因子选股模型更具有适用性。但考虑到实验数据获取的难易程度和论文实验时间的限制，本文将采取第二种方式对股票进行分类，并采用证监会的行业分类标准作为选取行业的参考。

二、行业代码和所属股票的选取

（一）研究行业的选取

作为拉动经济的三驾马车之一，消费在我国的经济中起着重要的作用，同时消费也和我们每个人的生活息息相关，所以消费一直就是资本市场中的“常青树”受广大投资者青睐，所以本文将选取消费作为研究的对象之一。但由于消费的范围太广，所以消费的概念比较笼统，因此本文将使用更为具体的泛消费行业作为实际的行业研究对象。在投资市场中，生物医药行业长期以来都被看作是朝阳行业。伴随着经济水平的不断提高，人们对于医药的需求越来越高，所以医药也同消费一样具有巨大的投资价值。近年来随着生物技术的突破，生物医药行业逐渐成为了主流，所以本文将生物医药行业作为研究的对象之一。

（二）行业代码的选取

选定了研究行业之后的下一步就是确定行业包含的股票，为了确定行业包含的股票，本文将采用证监会的二级行业分类标准来辅助确定。根据泛消费行业和生物医药行业的特点以及参考前人的研究^[12]，本文将证监会行业分类中 C13、C14、C15、C16、C17、C18、C19、C20、C21、C22、C23、C24、F51、F52、P82、H61、H62 等二级行业分类代码所属于的行业归为泛消费行业，将 C26、C27、Q38 等二级行业分类代码所属于的行业归为生物医药行业。其中泛消费行业和生物医药行业分类代码的具体说明详见附录 1。

（三）股票的选取

本文选取研究验证的时间范围是 2012 年 1 月 2 日至 2021 年 8 月 31 日，根据泛消费行业和生物医药行业所选取的分类代码，于 2012 年 1 月 2 日从我国 A 股中选取了所有这两个行业分类代码包含的股票作为初始股票池，其中泛消费行业一共选取了 279 只股票作为泛消费行业的初始股票池，生物医药行业一共选取了 202 只股票作为生物医药行业的初始股票池，其中两个行业所选取的股票详见附录 2。

第二节 候选因子选取

一、候选因子选取依据

股票市场是一个十分复杂的市场，在宏观上受到各种政策、国际重大事件的影响，在微观上还受到每个上市公司内部因素的影响。当前越来越多的研究者采用数量化的方式对股票市场进行研究，将各种对股市产生影响的因素通过可以量化的指标进行表示，以使得股票市场更加便利的使用统计学模型和各种各样的机器学习模

型进行分析。本文为了尽可能准确的刻画股票未来的收益，选取了通联量化数据平台上当前可以提供的所有因子，共计 232 个因子作为初始候选因子。将选取的因子按照因子的表达的含义进行分类，可以分为质量因子、情绪因子、风险因子、成长因子、价值因子、技术因子、动量因子等 7 个大类，其中选取的所有因子详见附录 3。

二、候选因子分类

（一）质量因子

质量因子通常是指那些刻画上市公司属性的因子，常见的质量因子有盈利能力、增长率、资本结构等。通常由上市公司多个不同维度的财务指标构成。

（二）情绪因子

投资者的情绪对于股票的走势会产生较大的影响，当投资者普遍看好某只股票的时候，往往这只股票的价格就会上涨。当投资者普遍不看好某只股票的时候，往往这只股票的价格就会下跌。因为投资者的情绪会影响到投资者的投资行为，进而影响到整个股票市场。因此本文为了尽可能掌握投资者情绪，把握股票未来的走势，故选取了情绪因子，其中常见的情绪因子有 TVSTD20、TVMA20、心理线指标等因子。

（三）风险因子

风险因子是一系列衡量股票风险指标的总称。衡量一个股票是否是具有投资价值的股票，需要充分考虑风险因子，从而尽可能的度量股票在未来一段时间内的不确定性以及系统风险。常见的风险因子有下跌波动、下跌贝塔等因子。

（四）成长因子

成长因子是一系列反应股票未来成长性的指标组成，现实中常常被用来区分股票的未来的收益。通常来说，将资金投资到成长性高的股票获得的收益会比将资金投入到低成长性的股票中多。常见的成长因子有净资产增长率、营业利润增长率等。

（五）价值因子

价值因子正如字面上的意思，是用于描述股票当前价值的因子，投资者可以通过价值因子去市场中寻找有投资潜力的股票。常见的价值因子有现金流市值比、收益市值比等因子。

（六）技术因子

技术因子是一类由技术指标转化而来的因子，和描述公司内部属性的质量因子相比，技术因子更注重的是分析市场行为，从市场行为中推断出股票价格变动的方向。常见的技术因子有 10 日指数移动均线、10 日移动均线等因子。

（七）动量因子

动量因子是将物理学中动量相关概念引入到了股票市场，认为股票的走势也具有动量，即当前股票是处于上涨状态，那么就认为股票未来仍会处于上涨状态。常

见的动量因子有八季度净利润变化趋势、10 日乖离率等因子。

大类因子包含部分二级因子如表 2-1 所示：

表 2-1 部分候选指标

大类名称	具体指标
质量因子	应付账款周转天数、应付账款周转率、对数总资产 账面杠杆等
情绪因子	心理线指标、相对强弱指标、10 日平均换手率等
风险因子	24 月累计收益、下跌贝塔、下跌波动、 收益相对波动等
成长因子	净资产增长率、归属母公司股东的净利润增长率、营业利润 增长率
价值因子	市盈率、市净率、市现率、分析师盈利预测等
技术因子	波幅中位数、K 日指数移动均线、资金流量指标等
动量因子	股票的 20、10、5 日收益、12 月相对强势、分析师推荐评 级变化、分析师盈收预测变化等

三、确定候选因子计算周期

为了获取到具体的因子值，需要对因子的计算时间进行规定，如果因子的计算时间明确的体现在因子名称中，那么就使用因子名称中的计算时间进行因子值的计算，比如 10 日均线、20 日均线等，明确的表明当日因子的值就是计算前 10 日、前

20 日收盘价的平均值。但是有一些因子的计算时间并没有在其因子名称中体现出来,那么我们就需要对这种因子的计算时间分情况规定,具体如下文所述:

1.当因子的计算时间有推荐的计算时间时,选取推荐计算时间中最小的值作为因子的计算时间。因为本文选取的因子数据是日度数据,为了数据尽可能的丰富,所以选取推荐时间范围内最小的时间作为因子的计算时间。比如心理线因子,可以计算 k ($k=1,2,3\cdots$) 日的心理线因子,但是推荐的心理线计算时间是 25 日到 75 日,那么此时就按照最小的 25 日作为心理线因子计算的时间。

2.当因子没有推荐的计算时间时,那么为了使得因子数据尽可能多样,就按照该因子最低可以计算的时间进行计算,比如对数市值因子(LCAP)的最低计算时间是 1 天,那么就采用 1 天作为计算时间,而下跌贝塔因子(Downside beta)最低计算时间是 12 个月,那么就按照 12 个月来计算。

第三节 数据预处理

数据预处理主要是对股票的原始数据和因子的原始数据进行处理,其主要目的是为了使得数据符合模型训练和测试的要求。数据预处理包含数据处理、因子选取和生成训练和测试数据集等三个环节。

一、数据处理

(一) 股票最终筛选

由于初始股票池中的股票是于 2012 年 1 月 2 日从 A 股票上获取的,虽然在 2012 年 1 月 2 日当天选取的股票是符合研究条件,但是这并不代表在整个选取的时间范围内这些股票也是符合研究条件,有些股票可能因为各种各样的原因出现异常甚至退市。所以需要对所选取的股票进行最终筛选,将所有出现过异常甚至退市的股票给剔除掉,其中股票的最终筛选主要分为如下两个步骤:

1.对两个行业股票池中的股票进行退市检查,查询股票池中的股票在本次研究选取时间范围内的最后一个交易日的上市情况,如果仍在 A 股中那么就将其保留。如果已经退市则将其剔除。

2.对上一步检查之后还剩下的股票,进行异常检查。在上一步中,可能有些股票在退市检查时股票是正常的,但是在所研究的时间范围内出现过异常,所以需要剩余的股票做进一步的异常检查。在研究的时间范围内,对每只股票每个交易日的状态进行检查,只要其出现过异常,就需要将这只股票剔除掉。

(二) 股票数据标签化

为了满足构建选股模型算法的训练和测试需要,需要将股票数据标签化处理。处理的方式主要分为如下三个步骤:

1. 计算股票每个自然月的月度涨幅。
2. 计算每个行业所有股票每个自然月的月度涨幅的平均值作为行业平均月度涨幅。
3. 将股票的月度涨幅和所属行业的平均月度涨幅进行比较，如果涨幅大于行业平均涨幅，就将股票当月的标签设置为 1，如果涨幅小于等于行业平均涨幅，就将股票当月的标签设置为 0。

股票月度涨幅的计算如公式（2-1）和（2-2）所示

$$DayRise_t = (Close_t - Close_{t-1}) / Close_{t-1} \quad (2-1)$$

其中 $DayRise_t$ 是 t 日涨幅， $Close_t$ 是 t 日收盘价， $close_{t-1}$ 是 $t-1$ 日收盘价。

$$MonthRise_k = \sum_{t=1}^{Month} DayRise_t \quad (2-2)$$

其中 $MonthRise_k$ 为第 k 个月度涨幅， $Month$ 自然月的天数， $DayRise_t$ 是 t 日涨幅

行业平均月度涨幅计算如公式（2-3）所示：

$$MonthRiseAerage_k = \frac{(\sum_{k=1}^N MonthRise_k)}{N} \quad (2-3)$$

其中 $MonthRiseAerage_k$ 是行业第 k 个月的平均涨幅， N 是行业包含股票的数量， $MonthRise_k$ 则是行业第 k 个股票的月度涨幅。

（三）数据标准化

股票经过了退市检查之后，剩下的股票都是正常的股票，所以不会因为股票退市导致股票和因子数据出现缺失值，同时股票的异常检查包括*ST 和 ST⁵检查，所以也不会由于股票出现*ST 和 ST 等状态导致股票和因子数据出现异常值。所以对于选取的股票和因子数据并不需要进行异常值和缺失值处理，只需进行标准化处理即可。

本文选取的因子数据由于各种因子的含义不同，导致不同因子之间的数据差别非常大，所以为了避免选股算法训练过程突出数值较高因子的作用，而削弱数值较

⁵ ST 和 *ST：当一个公司连续两年亏损或者净资产低于股票面值的时候，会在股票名字前面加上 ST（Special Treatment），用于警示投资者注意投资风险。当公司连续三年的经营都未有改善时，就会加上*，意为退市风险。

低因子的作用，保证结果的可靠性，需要对因子数据进行标准化处理。另外本文择时策略使用的强化学习算法中将股票五个维度的数据作为输入样本输入到神经网络中进行训练，但是其中的成交量和股票的价格之间数值差别太大，所以也存在上述问题，所以也需要对股票数据进行标准化处理。具体标准化的计算如公式（2-4）所示：

$$y = \frac{x - \min}{\max - \min} \quad (2-4)$$

其中，max 和 min 分别为因子样本上的最大值和最小值。

二、选取构建选股模型的因子

本文共计选取了 232 个因子作为初始候选因子，但是并不是每个因子都能够深刻的刻画股票未来的收益，有些因子用于反映股票未来收益的作用并不明显，即因子的值和股票未来的收益几乎无关，所以需要剔除这些因子。为了高效的筛选出用于构建选股模型的因子，本文主要分为了因子初筛和因子最终筛选两步来筛选因子，因子初筛和因子最终筛选具体方法和步骤为如下所述：

（一）因子初步筛选

对于当前因子值和未来收益之间的关系通常可以使用信息系数来判断，信息系数越大，那么该因子和股票未来收益之间的关系就越密切，所以本文将计算信息系数来初步筛选因子。本文选取的调仓⁶时间为一个月，所以计算信息系数中股票收益采用月度收益，同理因子也采用月度值进行计算，将信息系数值大于等于 10%和小于等于-10%的因子保留，将小于 10%和大于-10%的因子给剔除掉，其中信息系数的计算如公式（2-5）所示：

$$IC = Corr(A_n, P_{n+1}) \quad (2-5)$$

其中 IC 指的是信息系数， A_n 值的是第 n 期的因子月度值， P_{n+1} 指的是第 n+1 期的月度收益。

（二）因子最终筛选

因子经过初步筛选之后，留下的都是与股票未来收益之间关系较为密切的因子。为了进一步确定能够解释股票未来收益的因子，本文设计了一种新的因子选取方法作为因子最终筛选的方法。其中方法具体的步骤为如下五步所述：

第一步，首先分行业计算行业股票池中每个股票在训练周期⁷内对应因子的平均

⁶ 调仓：对所持有股票的仓位进行调整和股票品种的更换。

⁷ 训练周期：对于选股模型的训练，会将数据集切分为训练样本和测试样本，其中训练样本包含的时间跨度即为训练周期

值，然后以因子为单位生成股票对应因子值序列，最后在每个因子值序列中以因子值的大小对股票进行排序，其中股票的因子值序列如图 2-2 所示：



图 2-2 股票因子值序列

第二步，因子与股票未来收益之间的关系存在正向和反向两种关系。其中正向关系就是指因子值的越大股票未来的收益就越高，因子值越大的股票盈利的机会就越大，例如市值因子。而反向关系则是指因子的值越小股票未来的收益就越好，因子值越小的股票盈利机会就越大，例如应付账款周转率因子。所以根据因子对于股票未来收益正向关系的描述，对正向因子选取了因子值最大的 10 只股票作为因子的代表股票。根据因子对于股票未来收益反向关系的描述，对反向因子选取了因子值最小的 10 只股票作为因子的代表股票。

第三步，为因子的代表股票中的每一只股票分配 200 万初始资金，然后在训练周期的第一个交易日以当日开盘价和万分之五的手续费进行买入，并一直持有到训练周期的最后一个交易日。紧接着在最后一个交易日以当日的收盘价乘以当前持有股票的数量计算每只股票的市值。最后将每只股票当前的市值和手中剩余的资金求和，将得出的数字作为因子的得分。

第四步，计算行业平均得分。为行业内的每只股票分配 2000 万初始资金，然后在训练周期的第一个交易日以当日开盘价和万分之五的手续费进行买入，并一直持有到训练周期的最后一个交易日。紧接着在最后一个交易日以当日的收盘价乘以当前持有股票的数量计算每只股票的市值。最后将每只股票当前的市值和手中剩余的资金求和，将得到的数字在除以行业包含股票的数量，将得出的结果作为行业平均得分。

第五步，以因子为单位，将因子的得分和所属行业平均得分进行比较，如果大于行业平均得分，就认为这个因子是有较强的解释能力去解释股票未来的收益，就将其保留。如果小于等于行业平均得分，就认为这个因子没有解释股票未来收益的能力，就将其丢弃。因子最终筛选的流程如图 2-3 所示

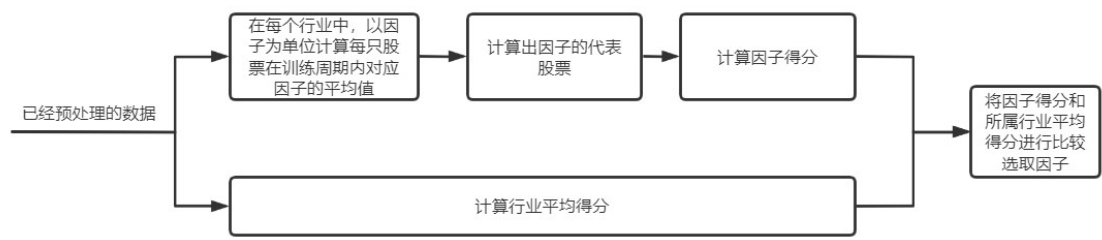


图 2-3 因子最终筛选流程

三、生成选股模型数据集

选取了构建选股模型的因子之后，将每个行业得到的因子作为最后选股模型训练的因子。然后将每个行业选取的因子分别构建每个行业的选股模型使用的数据集，其具体构造方法为下文所述：

- 1.首先获取每只股票获取所属行业选取的因子的标准化数据，并且以月为单位计算月每个因子的月平均值作为该因子的月度数据，从而每只股票形成了一条因子月度数据。
- 2.其次将股票的因子月度数据和股票下月的标签进行组合得到数据集的一条数据，具体的一条数据为图 2-4 所示：

第t月	股票1	(因子1的t月平均值, 因子2的t月平均值, 因子3的t月平均值, ..., 因子n的t月平均值, 股票1的t+1月标签)
	股票2	(因子1的t月平均值, 因子2的t月平均值, 因子3的t月平均值, ..., 因子n的t月平均值, 股票2的t+1月标签)
	股票3	(因子1的t月平均值, 因子2的t月平均值, 因子3的t月平均值, ..., 因子n的t月平均值, 股票3的t+1月标签)
	⋮	⋮
	股票k	(因子1的t月平均值, 因子2的t月平均值, 因子3的t月平均值, ..., 因子n的t月平均值, 股票k的t+1月标签)

图 2-4 t 月数据集

第四节 构建多因子选股模型

本文将采用当前在多因子选股中表现较好的 LightGBM、XGBoost 和 AdaBoost 三种算法分别构建不同的多因子选股模型，根据三种算法的表现，为不同行业的股票选择适合的算法构建多因子选股模型。

一、模型输入输出

LightGBM、XGBoost 和 AdaBoost 三种算法构建多因子选股模型的输入为上一步生成的数据集中将标签去除之后剩余的部分。

输出则要分几种不同的情况，当输出用于模型评估时输出为股票下个月的类别，当输出用于选股评估时，输出的则是下个月股票属于两种标签值的概率，因为在选股评估时，需要对模型选出股票的数量进行控制，如果选取股票的数量过多，那么大量股票频繁的交易所带来的手续费对最终的收益影响特别大。

模型在两种情况下的输入、输出如图 2-5 所示：

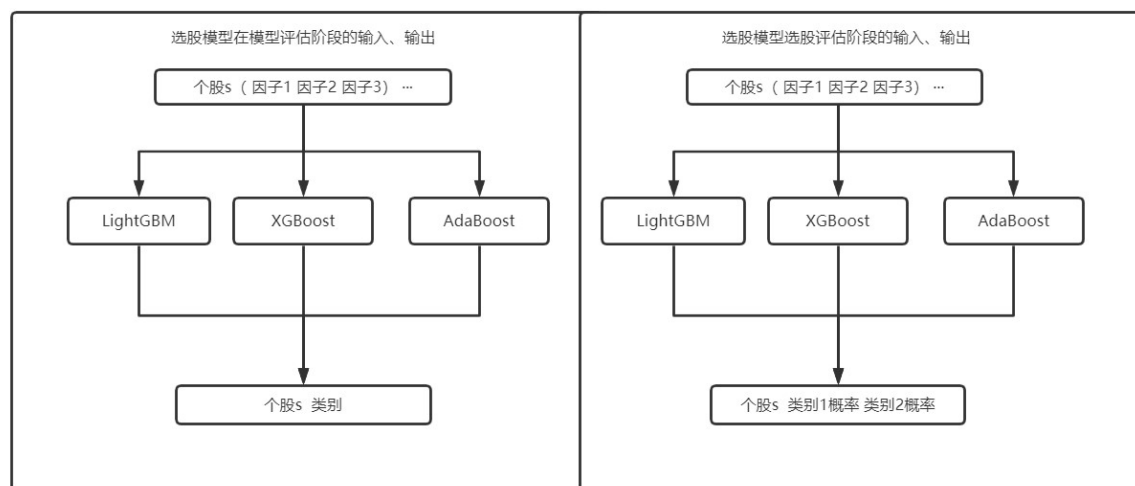


图 2-5 模型在模型评估和选股评估两种情况下的输入、输出

二、模型训练方式

本文采用滚动的方式进行训练，即本期测试的数据会出现在下一期的训练数据中。具体滚动的流程如下文所述：

1.滚动的窗口为 t ，即每一个训练样本由 t 个连续自然月的个股数据构成（序列的长度为 t ），个股的月数据包含经过最终筛选后的 n 个因子的月平均值，单个训练样本是一个 $t \times n$ 矩阵，输出则是 $t+1$ 个月的股票类别或者股票类别概率。所以训练周期内的总样本就是 k 个 $t \times n$ 的矩阵，其中 k 是行业包含的股票数量。

2.因为预测是以月度为基本时间单位进行,所以每次滚动的步长都是1,即窗口每次向后滑动1个月。所以滚动窗口的参数为 $size = t \times k$, $slide = 1$, 滚动训练具体的流程如图 2-6 所示:

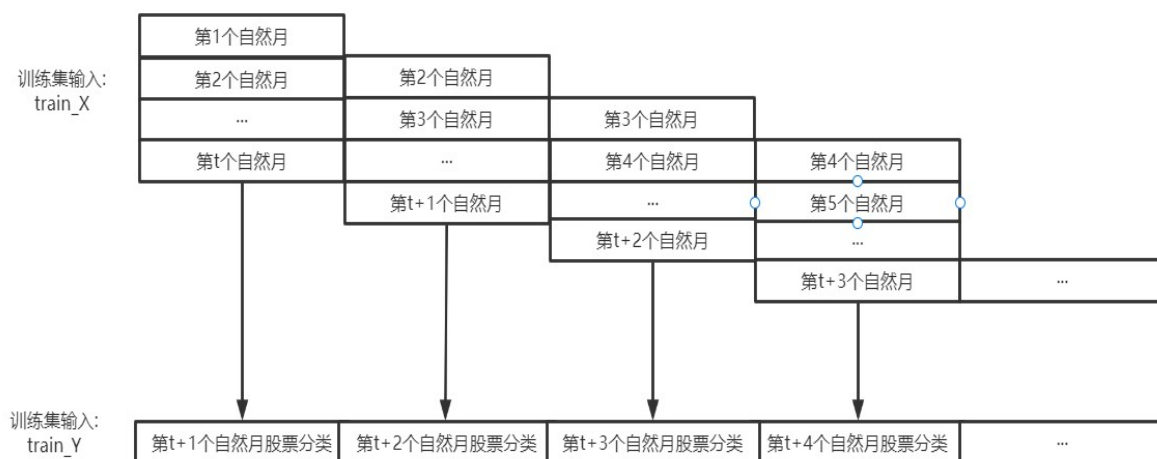


图 2-6 滚动训练流程图

三、模型训练过程

本文将生成数据集使用 LightGBM、XGBoost 和 AadBoost 三种算法进行训练,根据前人研究的基础上将 Adaboost 和 LightGBM 的学习率设置为 0.01,而 XGBoost 的学习率则设置为 0.05,并对三种算法独立进行训练。最后将每种算法预测的标签的概率进行输出,以 0.5 为概率类别分界值,当标签值为 1 的概率大于 0.5 的时候则为标签 1 的类别,否则为标签 0 的类别。

四、分类预测结果评估

(一) 模型评估

本文构建的多因子选股模型实例需要使用各种性能指标对其进行评估,以便查证选股模型实例的训练效果,针对模型的评估,本文采取了准确率和 AUC 两个指标,因为单独看准确率这个指标无法确定模型的真实优劣。假设一个极端例子,当样本有 9 个正例,1 个负例时,如果模型全部预测为正例,那么准确率也可以高达 90%,尤其是在这个负例样本为正例的概率还是最高的情况下,模型的性能本应极差,但从准确率上判断却是相反的结果。所以本文还引入了 AUC, AUC 可以很好的描述整体性能的高低。两个指标的具体含义为如下所述:

1.准确率 (Accuracy): 准确率是机器学习中常用的评价指标,它是一个用来评估机器学习算法识别能力的一个指标,准确率越高,那么算法识别的能力就越强,其计算公式如 (2-6) 所示:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (2-6)$$

其中 TP 表示分类正确的正例样本数， TN 表示分类正确的负例样本数， FP 表示分类错误的正例样本数， FN 表示分类错误的负例样本数。

2.ROC 曲线下面积（AUC）：这个指标是用来对模型识别能力的判断，特别是当样本并不均衡的时候，采用这个指标可以非常好的说明分类器的整体性能。

（二）选股收益评估

因为本文研究的对象是股票，而多因子选股的目的是为了取得较高的收益，所以本文不但需要对选股模型实例进行评估，同时还要对选股收益进行评估。对于选股收益的评估，本文采用的方法是将股票于 2020 年 2 月 3 日到 2021 年 8 月 31 日的测试区间内进行回测，但是不采取任何的择时策略，具体回测的步骤如下文所述：

1.每期输出的股票类别中，对股票以标签值为 1 的概率进行排序，选取排名前 10 的股票作为本期多因子选股模型选取的股票。当有概率相同的时候，就对股票计算训练周期内的平均成交量，选取平均成交量较大的股票作为本期选股。成交量大的股票流通性更高，流通性高的股票在真实交易环境中买入和卖出时出现滑点⁸的概率相对较小，使得做出的交易最大程度的符合预期。

2.设置股票的初始可用资金为 2000 万，然后在回测区间的第一个交易日为第一期选取的股票进行初始资金分配，以平均分配的方式进行分配，但是每只股票最后分配的资金都是整数，其具体如公式（2-7）所示：

$$p = \lfloor s/n \rfloor \quad (2-7)$$

$$s = s - p * n \quad (2-8)$$

其解释如下： s 代表当前手中的资金，一开始就是初始资金 2000 万，而 n 代表本期的股票总数， p 代表单个股票的分配的向下取整资金。在经过一次分配之后手中剩余的资金即为公式 2-8 所计算之后的结果。

3.如果当前交易日是回测区间上的第一个交易日，那么以当前交易日的开盘价对股票进行买入，为了尽可能的模拟真实交易市场，为买入股票设置一个万分之五的手续费。

4.如果当前交易日是调仓日，那么首先获取在上一期存在但是本期不存在的股票。然后对这些股票进行判断，如果当前交易日这些股票是处于持有状态，那么就

⁸ 滑点：是指股票买入和卖出时预期的交易价格和实际的交易价格存在较大的差别的一种现象

将这些股票强制卖出，并且收回其资金到剩余资金 s 中。为了尽可能的模拟真实交易市场，在卖出股票的时候设置一个万分之五的手续费。紧接着对本期存在但是上一期没有存在的股票进行资金的初始分配，分配规则同样是平均分配。最后将它们以本交易日的开盘价进行买入。

5. 每一个交易日都将每只股票当日的收盘价乘以该股票的持仓数量计算出当前股票的市值，然后将所有股票的市值与当前手头剩余的资金 s 求和，最后将计算的总和作为当前交易日的收益。

为了对选股收益进行评估，本文选取了年化收益率和最大回撤等两个指标对获取的收益进行评估，其中两个指标的计算方式如下：

1. 年化收益率

年化收益率（Annualized Returns），这个指标是以年为单位进行计算，如果不满一年需要将其换算成年。指标表示投资期限为一年时，从年初到年末收益的变化情况，其具体计算如公式（2-9）所示：

$$AR = \left(\frac{V}{I}\right)^{\frac{t}{y}} - 1 \quad (2-9)$$

V 为投资时间 t 后的市值， I 为投资本金， y 为一年的有效投资时间，对于股票期货市场 $y=250$ 日。

2. 最大回撤率

最大回撤率（Max Drawdown），描述买入投资组合后可能出现的最糟糕的情况。

$$MaxDrawdown_t = \max\{ (D_i - D_j) / D_i \} \quad (2-10)$$

D_i 为策略组合第 i 天的净值， D_j 策略组合第 j 天的净值， T 为回测交易日数量； $T \geq j \geq i$ 。

第五节 构建择时策略

使用多因子选股模型进行选股的目的是期望将未来一段时间内整体收益较好的股票挑选出来，但是选出的股票仍然会存在许多涨跌波动。所以会出现买入一个整

体收益较好的股票，却因为交易时机不对导致最终获取的收益不理想，甚至变成负收益的情况。所以为了尽可能的避免出现这种情况，本文针对行业的特点设计了一种基于长期均线和深度强化学习的择时策略，并对分行业构建的选股模型选取的股票进行回测。

传统技术指标由于自身的局限性，导致基于传统指标构建的择时策略在实际的交易中并不理想。比如常见的均线策略，通过使用长期均线来判断股票未来半年或者一年的综合趋势⁹：上涨、下跌或者震荡。然后在上涨的趋势中使用短期均线和中期均线来避免短期之内的回调¹⁰，从而提高股票交易的收益。但是由于均线本身计算的原因，所以均线发出的信号通常来说都是晚于真实股票的走势。对于长期均线发出的信号来说，这种滞后的消息是可以接受的，但是对于短期和中期均线发出的信号来说，这种滞后往往不能够接受，因为这种滞后有可能造成期待的买入和实际的买入是截然相反的情况。

近年来，将深度强化学习用于构建择时策略的研究也逐渐变多，从董焕彬^[10]和于龙飞^[11]等人的论文中可以看到针对短期的股票市场进行学习，深度强化学习构建的择时策略可以取得较好的结果。但是股票市场过于复杂，长时间的股票市场中包含的金融噪声成倍提高，导致深度强化学习使用长时间的数据学习构建的择时策略取得的效果并没有使用短期数据情况下取得的好。

针对基于均线的择时策略和基于深度强化学习的择时策略出现的问题，本文提出了一种解决方案，将均线择时策略的优点和深度强化学习择时策略的优点进行结合，设计了一种新型择时策略。并根据不同行业的特点，对择时策略进行了优化，以使其用于分行业构建的多因子选股模型选取的股票中。

一、构建长期均线策略

（一）长期均线

长期均线一般指的是 120 到 240 日，每年的交易日大约有 240 天左右，所以取 120 到 240 日中的任意数计算均线值，并根据均线值和当日收盘价的上穿和下穿的关系可以大致预测未来半年到一年的时间内股票整体是上升还是下跌趋势，其具体判断是上升和还是下跌趋势为如下描述：

1. 当日收盘价大于当日均线值并且昨日收盘价小于昨日均线值，预测股票未来是上升的趋势。
2. 当日收盘价小于当日均线值并且昨日收盘价大于昨日均线值，预测股票未来是下跌的趋势。

⁹ 综合趋势：即股票一段时间的走势，从这段时间的起点到终点所连的线的斜率为正，则是上升趋势，为负则是下跌趋势，其中一段时间通常指半年以上的时间。

¹⁰ 回调：价格在上涨趋势中，价格上涨速度太快，受到卖方打压而暂时回落的现象。

（二）构建长期均线策略

本文通过设置均线值为 120 到 240 日之间的每一个数，然后使用收盘价和均线值之间的关系发出买入和卖出信号，通过对需要择时月份前七年的历史数据采用这个信号进行回测，最终将获得的收益作为均线值的得分，选取得分最高的均线值作为择时月份所使用的长期均线。而收盘价和均线之间产生发出的信号则如下文所定义：

1. 当日收盘价大于当日均线值并且昨日收盘价小于昨日均线值，发出买入信号。
2. 当日收盘价小于当日均线值并且昨日收盘价大于昨日均线值，发出卖出信号。

二、构建深度强化学习 DQN 择时策略

（一）定义深度强化的环境、动作和奖励

本文在中短期择时策略采用深度强化学习 DQN 算法用于构建择时策略，股票市场不同于往常的游戏环境，对于智能体¹¹所需的三要素：环境、可选的动作、奖励没有标准化的答案，对三者的定义即使有微小的差别，最后求出的结果都可能会截然不同。所以如何构建 DQN 算法智能体所处的环境、可执行的动作以及执行动作之后的奖励，对于训练一个优秀的智能体就显得十分重要。为了尽可能的使得智能体掌握股票市场，本文对环境、可选动作和奖励作了如下定义：

1. 环境（Environment）：由于股票市场是一个复杂的市场，股票每天的涨跌并不是随机的，与之前的股票之前的状态有一定的关联。所以本文为了尽可能的描述股票历史状态和当前状态之间的关系，采用了股票的开盘价、收盘价、最高价、最低价、和成交量等五个维度的数据作为股票市场的环境。

2. 可选的动作（Action）：在现实中，对股票的交易可以分为买入、卖出和观望等三种动作，所以本文为了贴近现实，也将采取这三种动作作为智能体的可选动作。

3. 奖励（Reward）作为智能体执行动作之后的反馈，同样本着贴近现实的理念，采取股票的涨跌幅作为奖励。当智能体执行买入动作时，如果另一天股票是上涨，那么奖励就按照正常的涨跌计算。当智能体执行卖出动作时，那么奖励需要按照涨跌幅乘以-1 进行计算，即另一天下跌时，智能体的奖励为正，另一天上涨时，智能体的奖励为负。

（二）模型的输入输出

如上描述的那样，本文使用股票每日五个维度的数据作为当天的状态（State）输入到深度神经网络中进行训练，通过深度神经网络来解决 Q 值表爆炸的问题，以近似估计的方式计算每个动作的 Q 值。具体的输入和输出描述如图 2-7 所示：

¹¹ 智能体：强化学习中行为的决定者。

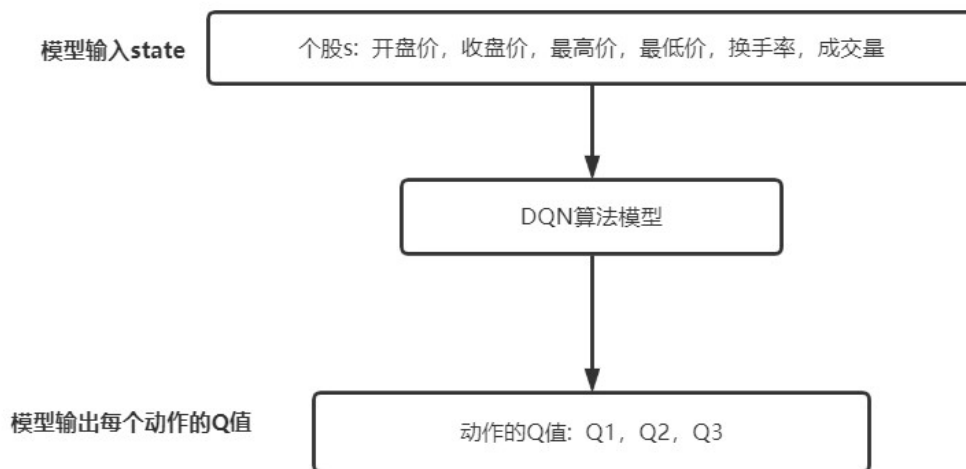


图 2-7 DQN 算法的输入和输出

（三）深度强化学习 DQN 模型的训练

本文采用滚动的方式进行训练，以月度作为滚动的单位，当前月份择时所使用的模型采用当前月份前一段时间的数据进行训练。而到了下一个月份，则本月则会作为训练数据用于训练，其中前一段时间是根据行业的特点来确定。另外对于 DQN 神经网络的参数设置中将 Gamma 值设置为 0.92，学习率则设置成了 0.01，而经验池的大小则是综合考虑实验设备内存大小设置以及数据集的大小设置成了 1000。

三、构建融合策略

本文采用长期均线和深度强化学习 DQN 算法来共同构成择时策略，以长期均线来判断上升还是下跌的趋势，然后在上升趋势中使用深度强化学习 DQN 算法来进行短期择时，其主要分为如下两种情况：

1. 当日收盘价大于当日的均线值，此时发出未来一段时间内将上涨的信号，则当前交易日深度强化学习 DQN 发出的信号都是有效的。
2. 当日收盘价小于等于当日的均线值，此时发出未来一段时间内将下跌的信号，则当前交易日深度强化学习发出的买入和持有信号是无效的，如果当前交易日还持有股票则在下一个交易日以开盘价进行强制卖出。

四、择时策略评估

将分行业构建的择时策略应用于分行业构建多因子选股模型选取的股票上，其操作方法和选股评估阶段的操作方法类似，但是本次需要使用构建的择时策略进行择时，其中具体方法为下文所述：

1. 首先股票设置股票的初始持有资金为 2000 万，然后在回测区间的第一个交易日为第一期选取的股票进行初始资金分配，以平均分配的方式进行分配，但是每只股票最后分配的资金都是整数，其具体的公式描述如下：

$$f = \lfloor m/n \rfloor \quad (2-11)$$

$$m = m - f * n \quad (2-12)$$

其解释如下： m 代表当前手中的资金，一开始就是初始资金 2000 万，而 n 代表本期的股票总数， f 代表带个股票的分配的资金。在经过一次分配之后手中剩余的资金即为公式 2-12 所计算之后的结果。

2. 如果当前交易日是调仓日，那么首先获取在上一期股票存在但是本期中不存在的股票。然后对这些股票进行判断，如果当前交易日是持有这些股票的话，那么就将这些股票强制卖出，并且收回其资金。为了尽可能的模拟真实交易市场，在卖出股票的时候设置一个万分之五的手续费。紧接着对本期存在但是上一期没有存在的股票进行资金的初始分配，分配规则同样是平均分配。

3. 在每一个交易日以长期均线和深度强化学习构建的择时交易信号进行交易。在交易完成以后以当日的收盘价乘以股票的持仓数量计算出当前股票的市值，并累加所有股票的市值与当前手头剩余的资金，最后将计算的总数作为当前交易日的收益。

对于择时策略收益的评估，本文不仅选取了选股阶段所使用的评价指标，还另外新增了年化超额收益率和夏普比率两个新指标。

1. 年化超额收益率

年化超额收益率（Excess of Market Average Return），代表通过主动投资所获得超过基准收益的部分，其中超额收益率有两种计算方式一种是绝对超额收益率，它通过将策略带来的收益直接减去基准收益的方式计算而来。另一种是相对超额收益率，它更加具有现实价值，其计算公式如（2-13）所示：

$$EMAR = \frac{(1 + AR)}{(1 + MR)} - 1 \quad (2-13)$$

其中 AR 为年化收益率， MR 为基准年收益率。

2. 夏普比率

夏普比率（Sharpe Ratio），这个指标是用来形容风险和收益之间的关系，即表示每承受一单位的风险可以获得的收益，夏普比率越大表示在相同的风险下所获得的收益越高。

$$SharpeRatio = (E(R) - R_f) / \sigma \quad (2-14)$$

其中 $E(R)$ 为策略组合的预期报酬率， R_f 为无风险利率， σ 为策略组合的标准差。

第三章 选股模型和择时策略实证

本章主要对第二章设计的内容进行实证，主要可以分为如下四个部分：

第一个部分，实证对象的获取和实验环境的确定。这个部分主要是获取实验数据与实验环境的确定。

第二个部分，数据预处理。这个部分主要是利用第二章的数据预处理中的步骤和方法对数据进行处理，其中经过异常检查后两个行业还剩下 451 只股票可以进行实验。而股票数据标签化则一共生成了 51865 个标签值。对因子的筛选也是按照第二章所规定的方法进行筛选，最后总共为泛消费行业和生物医药行业选取了 12 个和 16 个因子用于构建多因子选股模型。同时为了对比分行业构建多因子选股模型的效果，本文为此还新增一个特殊行业：不分行业，即将泛消费行业和生物医药的行业的数据合并，并以构建泛消费行业和生物医药行业多因子选股模型的方式构建不分行业的多因子选股模型，其中不分行业经过因子最终筛选选取了 12 个因子。

第三个部分，选股模型的训练和对比。这个部分主要是对第二章中设计的选股模型进行实证和结果对比。

第四个部分，构建择时策略和回测验证。这个部分主要也是根据第二章设计择时策略进行实证训练和结果对比。

第一节 实证对象获取和实验环境确定

一、数据获取

实证数据中的股票数据通过聚宽量化交易平台（<https://www.joinquant.com>）获取，而实证数据中的因子数据则是通过通联量化交易平台（<https://uqerdatayes.com/>）的数据 API 进行获取。

根据第二章中确定的行业以及行业包含的证监会二级分类代码，本文选取了 281 只股票作为泛消费行业的初始股票池，选取了 202 只股票作为生物医药行业的初始股票池。所本文针对这 481 只股票，逐个的从聚宽量化交易平台上获取了 2012 年 1 月 4 日到 2021 年 8 月 31 日期间的包含开盘价、收盘价、最高价、最低价和成交量等五个维度的日数据。为了降低偶然性对于本文研究的影响，本文选取了较广的时间范围，选取的时间内我国 A 股市场发生了许多重要事件，所以几乎包含了所有的行情趋势。同时从通联量化交易平台上逐个的获取了每只股票 232 个因子于 2012 年 1 月 4 日到 2021 年 8 月 31 日期间的日数据。

二、确定实验环境

由于 Python 拥有众多的库，十分易用，所以本文选股模型和择时模型均采用

Python 来实现，机器学习框架则使用 TensorFlow2.0，而回测¹²系统则考虑设计的复杂性以及运算的速度两方面，本文采用 Java 来构建本文的回测系统。

第二节 数据预处理

一、数据处理

（一）股票数据最终筛选

根据第二章中所描述的方法，首先对选取的股票在 2021 年 8 月 31 日检查所有股票是否仍在 A 股中，将不在 A 中的股票剔除。然后将剩余股票于 2012 年 1 月 4 到 2021 年 8 月 31 日期间的每个交易日检查是否出现 ST 和 *ST¹³，将期间出现过 ST 和 *ST 的股票进行剔除。经过处理之后泛消费行业一共保留了 263 只股票，而生物医药行业则保留了 188 只股票。

（二）股票数据标签化

计算泛消费行业和生物医药行业每只股票的月度涨幅，同时计算泛消费行业和生物医药行业每个月的行业平均涨幅，将泛消费行业和生物医药行业每只股票的月平均涨幅和所属行业的月平均涨幅进行比较，当股票的月度涨幅大于所属行业的月度平均涨幅时，则标签值设置为 1，否则标签值设置为 0。经过计算与比较之后，一共生成了标签值为 1 的数据 25391 条。一共生成了标签值为 0 的数据 26474 条。

（三）数据标准化

按照第二章中的标准化方法，对股票数据和因子数据进行标准化处理，其中标准化的主要分为三个步骤。第一步，求得数据的平均值 m 。第二步，求得数据的最大值 \max 和最小值 \min 。第三步，将每个数据减去平均值之后除以最大值减去最小值。

二、因子选择

（一）因子初步筛选

根据第二章因子初步筛选的方法，以股票为单位对每个因子在 2012 年 1 月至 2021 年 8 月共计 116 个月上计算信息系数，如果因子有一个月的信息系数大于 -10% 或者小于 10% 则剔除这个因子。将每个行业下的所有股票保留的因子取交集，将取交集之后剩余的因子作为显著性因子保留，其中泛消费行业一共获取了 117 个显著性因子，生物医药行业则获取了 57 个显著性因子，不分行业获取了 79 个显著性因子。

（二）因子最终筛选

首先根据第二章描述的因子选取方法，分别计算泛消费行业、生物医药行业和不分行业三者每个因子的代表股票，其次将因子的代表股票每只分配 200 万资金，于 2012 年 1 月 4 日以开盘价进行买入，一直持有到 2019 年 12 月 31 日，并以 2019 年

¹² 回测：根据历史数据来验证交易策略的可行性和有效性的过程

¹³ ST 和 *ST：当一个公司连续两年亏损或者净资产低于股票面值的时候，会在股票名字前面加上 ST（Special Treatment），用于警示投资者注意投资风险。当公司连续三年的经营都未有改善时，就会加上*，意为退市风险。

12月31日的收盘价分别乘以每只股票的持有数量计算出因子的代表股票的总市值作为因子的得分。然后对行业下每只股票以2000万资金于2012年1月4日以开盘价买入,并以2019年12月31日的收盘价乘以每只股票持有的数量计算出每只股票的市值,将所有股票的市值进行累加并除以行业包含股票的数量,将得到的数字作为行业的平均得分。经过计算得泛消费行业的平均得分是31605953.93、生物医药行业的平均得分为36296348.22,不分行业的平均得分为33951151.08。最后以行业为单位,分别将因子的得分与所属行业的平均得分进行比较,最终泛消费行业一共选取了12个因子用于构建多因子选股模型,生物医药行业则选取16个因子用于构建多因子选股模型,不分行业选取了12个因子用于构建多因子选股模型,其中每个行业选取的因子如表3-1至3-3所示:

表 3-1 泛消费行业选取的因子汇总

因子名称	因子类别
DAVOL20 20日平均换手率与120日平均换手率	情绪因子
InvestCashGrowRate 投资活动产生的现金流量净额增长率	成长因子
MA5 5日移动均线	技术因子
ROE 权益回报率	质量因子
BackwardADJ 股价向后复权因子	风险因子
ATR6 6日均幅指标	情绪因子
CMO 钱德动量摆动指标	动量因子
VDEA 计算VMACD因子的中间变量	情绪因子
PVT6 因子PVT的6日均值	动量因子
TRIX5 5日收盘价三重指数平滑移动平均指标	动量因子
Ulcer10	技术因子
JDQS2 阶段强势指标	情绪因子

表 3-2 生物医药行业选取的因子汇总

因子名称	因子类别
HBETA 历史贝塔	风险因子
MA5 5日移动均线	技术因子
MAWVAD	情绪因子

因子 WVAD 的 6 日均值	
OperCashGrowRate 经营活动产生的现金流量净额增长率	成长因子
SUOI 未预期毛利	成长因子
ILLIQUIDITY 收益相对金额比	技术因子
ATR6 6 日均幅指标	情绪因子
BIAS60 60 日乖离率	动量因子
KDJ_D 随机指标	技术因子
STM 计算 ADTM 因子的中间变量	情绪因子
SRMI 修正动量指标	动量因子
CMO 钱德动量摆动指标	动量因子
VEMA10 成交量的 10 日指数移动平均	情绪因子
VOSC 成交量震荡	情绪因子
minusDI 下降指标	情绪因子
JDQS2 阶段强势指标	成长因子

表 3-3 不分行业选取的因子汇总

因子名称	因子类别
DAVOL20 20 日平均换手率与 120 日平均换手率	情绪因子
InvestCashGrowRate 投资活动产生的现金流量净额增长率	成长因子
MA5 5 日移动均线	技术因子
NonCurrentAssetsRatio 非流动资产比率	质量因子
ROE 权益回报率	质量因子
ATR6 6 日均幅指标	情绪因子
BIAS60 60 日乖离率	动量因子
KDJ_D 随机指标	技术因子
SRMI 修正动量指标	动量因子
CMO 钱德动量摆动指标	动量因子

minusDI	技术因子
JDQS2 下降指标	成长因子

第三节 选股模型训练和对比

一、LightGBM、XGBoost 和 AdaBoost 模型训练

（一）数据集划分

为了满足训练使用的数据集大小，本文将 2012 年 1 月 4 日至 2019 年 12 月 31 日的日期划分为初始训练集，而 2020 年 1 月 2 日到 2021 年 8 月 31 日期间的数据作为测试集。并且为了动态更新模型的参数，本文以月度为单位对模型进行滚动训练，即初始训练数据是 2012 年 1 月 4 日到 2019 年 12 月 31 日，而第二期训练的数据则是从 2012 年 2 月 1 日到 2020 年 1 月 31 日期间的数据。其中对模型进行滚动训练如图 3-1 所示



图 3-1 数据集滚动划分图

（二）LightGBM、XGBoost 和 AdaBoost 多因子选股模型训练

LightGBM、XGBoost 和 AdaBoost 三个都是集成学习算法，在近年来的多因子选股研究中使用这三种算法均取得了较为不错的收益，所以本文选取这三种算法作为构建多因子选股模型的算法。将泛消费行业、生物医药行业和不分行业分别使用三种算法进行训练，其中泛消费行业训练模型所使用的特征为泛消费行业选取的 12 个因子，生物医药行业训练的特征则是生物医药行业选取的 16 个因子，不分行业训练的特征则是不分行业选取的 12 个因子。其中所有算法在测试集上的输出均为股票下个月的分类。利用训练好的模型进行选股时则是输出股票下个月属于每种标签的概率。训练的具体流程如图 3-2 所示

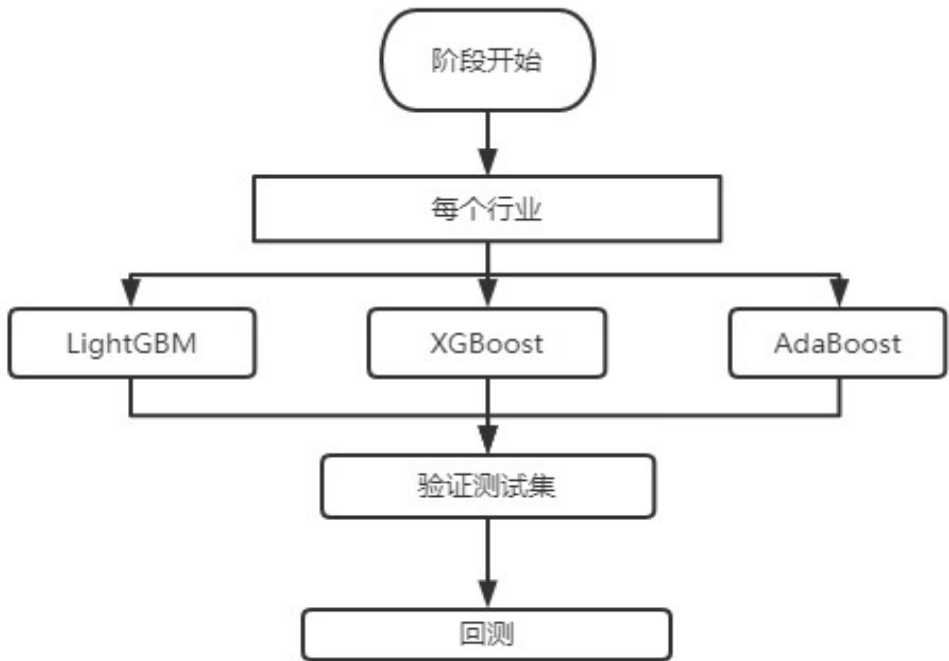


图 3-2 泛消费行业生物医药行业和不分行业三种进行模型训练流程图

本文以滚动的方式对每个月股票的分类进行预测，其中滚动的时间窗口为 $t=96$ ，即用前 96 个月的股票因子时间序列预测第 97 个月股票的类别，具体如图 3-3 所示：

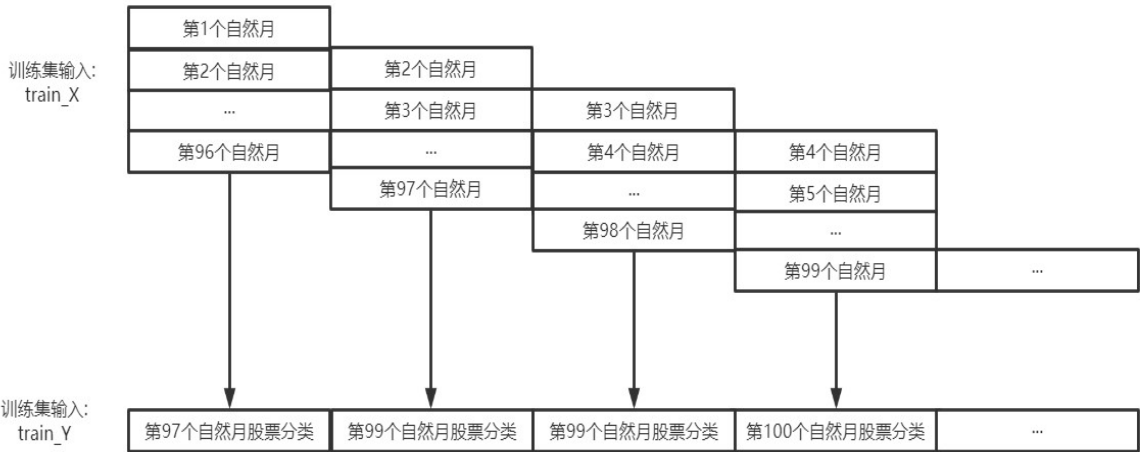


图 3-3 模型滚动训练图

二、训练过程与预测效果对比

本文对泛消费行业、生物医药行业在 LightGBM、XGBoost 和 AdaBoost 三种算

法下构建股票的分类预测模型。为了与泛消费行业和生物医药行业的结果进行对比，同时也对不分行业以相同的方式在三种机器学习算法下构建选股模型。另外本文还将从模型评估和选股评估两个角度对多因子选股模型进行评价，其中模型评估的目的是为了测试三种算法的训练的效果，而选股评估则是为了验证基于行业构建多因子选股模型的有效性。

（一）模型评估

本文将三个行业分别对三种不同的机器学习算法进行训练，19个月的滚动训练三个行业、三种算法一共得到了171个多因子选股实例股模型，由于实例模型过多，所以只选取最后一次滚动训练的模型进行评估。通过使用第二章中选取的模型评估指标对最后一期训练的模型进行相应指标的计算，其中最后一期模型两个指标计算的结果如图3-4到3-5所示：

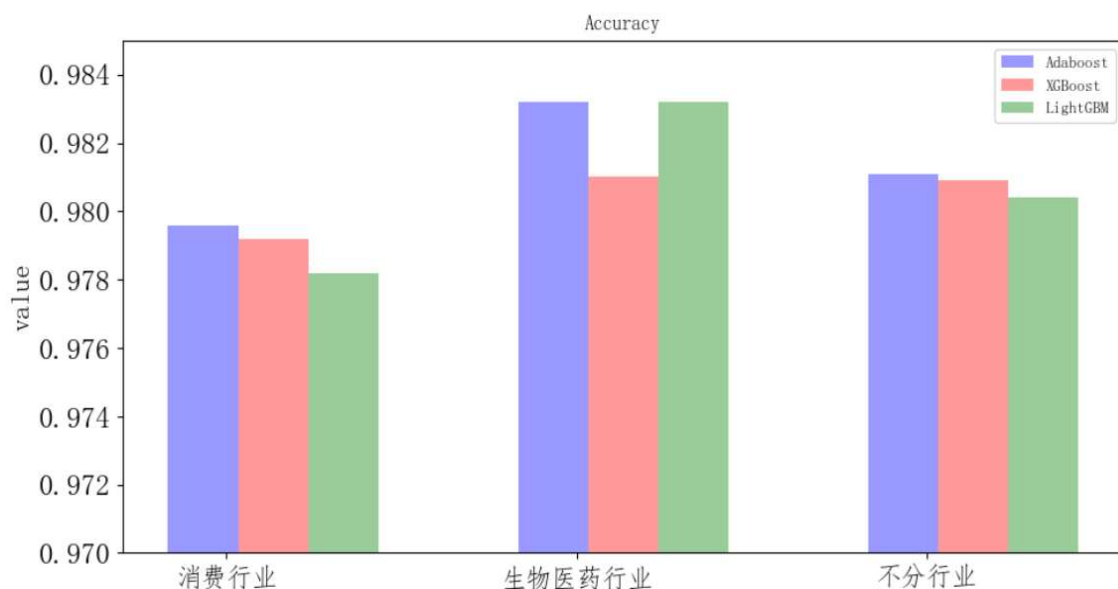


图 3-4 Accuracy

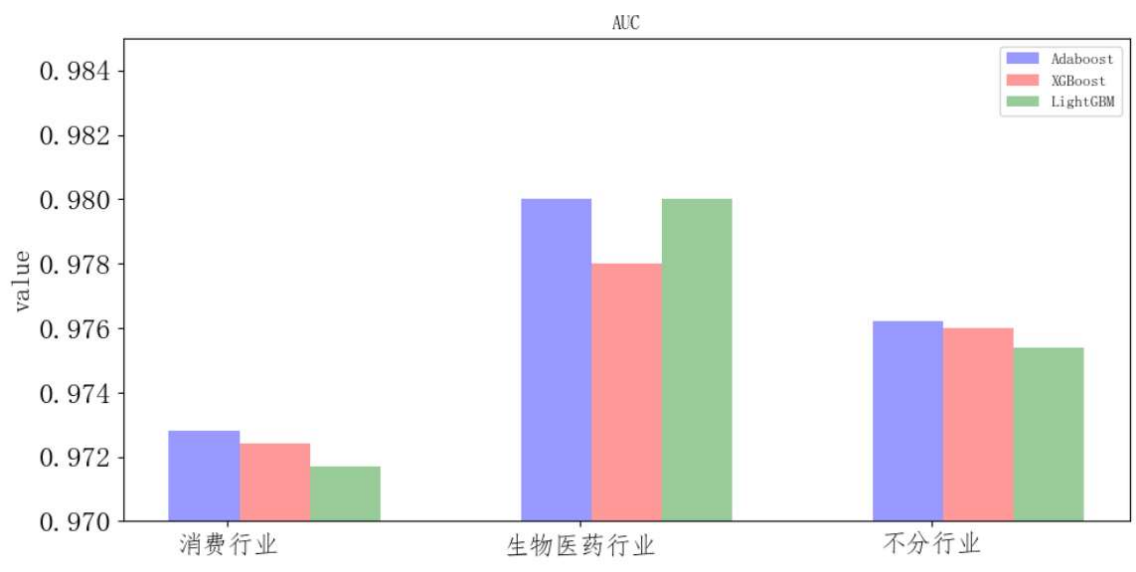


图 3-5 AUC

通过上图我们可以看到无论是哪种行业，得到的准确率（Accuracy）都是在 0.97 至 0.99 之间，准确率比较高。从图 3-5 中可以看到 AUC 的指标的值也是居于 0.97 至 0.99 之间，AUC 值也比较高。结合准确率和 AUC 两个指标的值，可知模型的整体分类性能比较出色，说明这些算法很好的找到了因子与股票未来收益之间的关系，可以将这些算法用于构建基于行业的多因子选股模型。

（二）选股收益评估

通过第二章选股收益评估的方法，对三个行业的收益进行计算得到了每个行业三种模型收益对比图，具体为图 3-6 至 3-8 所示：



图 3-6 泛消费行业三种模型选股收益对比图

结果图 3-6 中可以清楚看到 XGBoost 算法构建的多因子选股模型所选取的股票获取的收益都远远高于 LightGBM 和 AdaBoost 算法构建的多因子选股模型所选取股票获取的收益。由此可知,对于泛消费行业来说,使用 XGBoost 对其进行构建多因子选股模型是最合适的算法,所以本文将采用 XGBoost 算法选取的股票作为泛消费行业多因子选股模型选取的股票,并将其于其他两个行业的收益进行对比。

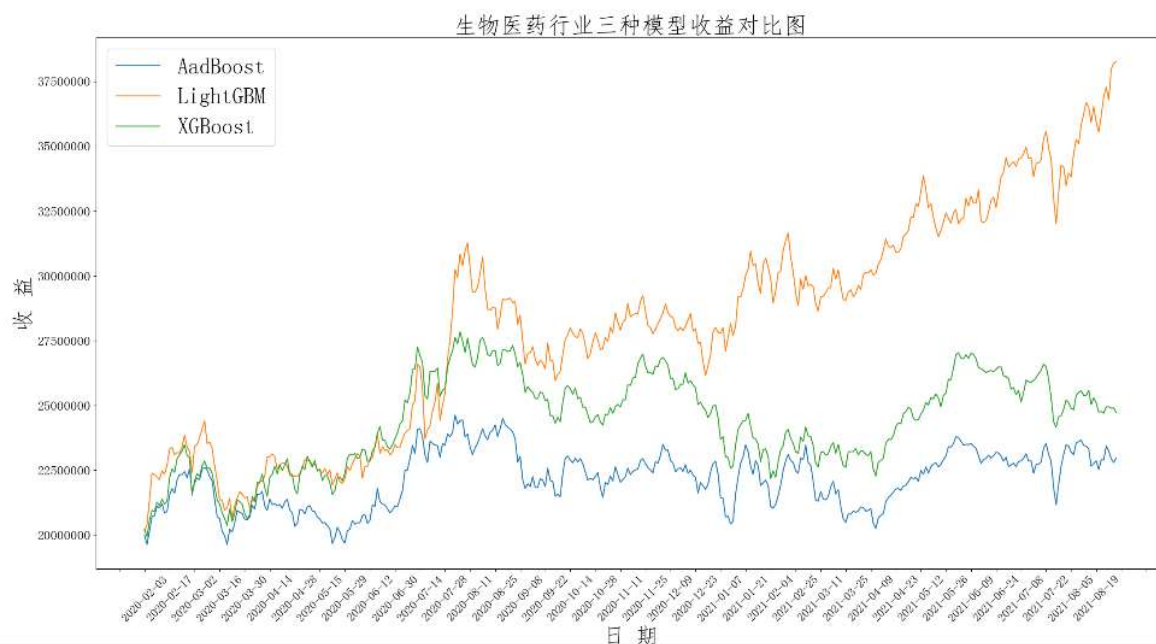


图 3-7 生物医药行业三种模型选股收益对比图

从图 3-7 中可以清楚看到 LightGBM 算法构建的多因子选股模型所选取的股票收益在 2020 年 2 月至 2020 年 8 月之间与 XGBoost 算法选取的股票获取的收益几乎一致,但均要好于 AdaBoost 算法选取的股票获取的收益。但在 2020 年 8 月份至 2021 年 8 月份的时间内,LightGBM 选股模型的收益比其他两种选股模型的收益都高,并且是越来越高。由此可知,对于生物医药行业来说,使用 LightGBM 算法对其进行构建多因子选股模型是最合适的算法,所以本文将采用 LightGBM 算法选取的股票作为生物医药行业多因子选股模型选取的股票,并将其与其他行业进行对比。

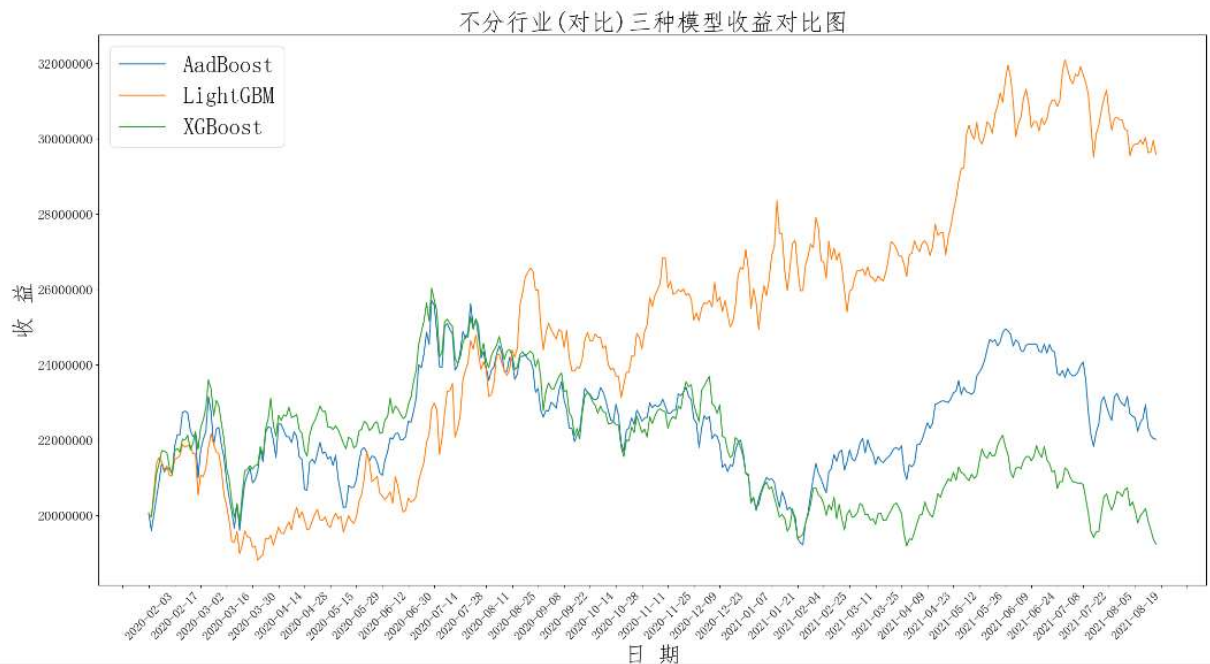


图 3-8 不分行业（对比）三种算法选股收益对比图

从图 3-8 中可以清楚看到 LightGBM 算法构建的多因子选股模型所选取的股票的收益在 2020 年 2 月至 2020 年 8 月期间比 XGBoost 和 AdaBoost 算法构建的选股模型选取的股票的收益都要低，但从 2020 年 8 月开始，LightGBM 选股模型的收益开始反超其他两种模型的收益，并且是在其他两种模型呈现下跌的情况，LightGBM 选股模型的收益仍保持增长。虽然 LightGBM 算法在后期表现较为优异，但是单单从图中仍然无法确定哪种算法更适合不分行业构建多因子选股模型。所以本文对三种模型取得的收益计算了年化收益率和最大回撤两个指标的数据，具体为表 3-4 所示：

表 3-4 不分行业三种算法收益指标

	收益指标	
	年化收益	最大回撤
AdaBoost	6.4%	25.22%
XGBoost	-2.58%	26.28%
LightGBM	29.64%	15.13%

根据表 3-4 中的结果可知，LightGBM 算法的年化收益是三种模型中最高的，而最大回撤却是三种模型中最低的。综合图 3-8 和表 3-4 的结果可知，使用 LightGBM 算法构建多因子选股模型是不分行业最合适的算法，所以本文将采用 LightGBM 算法选取的股票作为不分行业多因子选股模型选取的股票，并将其与其他行业进行对比。

将泛消费行业多因子选股模型选取的股票、生物医药行业多因子选股模型选取的股票和不分行业多因子选股模型选取的股票的收益进行对比，同时为了验证是否分行业构建多因子选股模型选股收益有实际意义，本文还加入了两个行业股票的平均收益作为基准收益进行对比，具体对比结果为图 3-9 和表 3-5 所示：

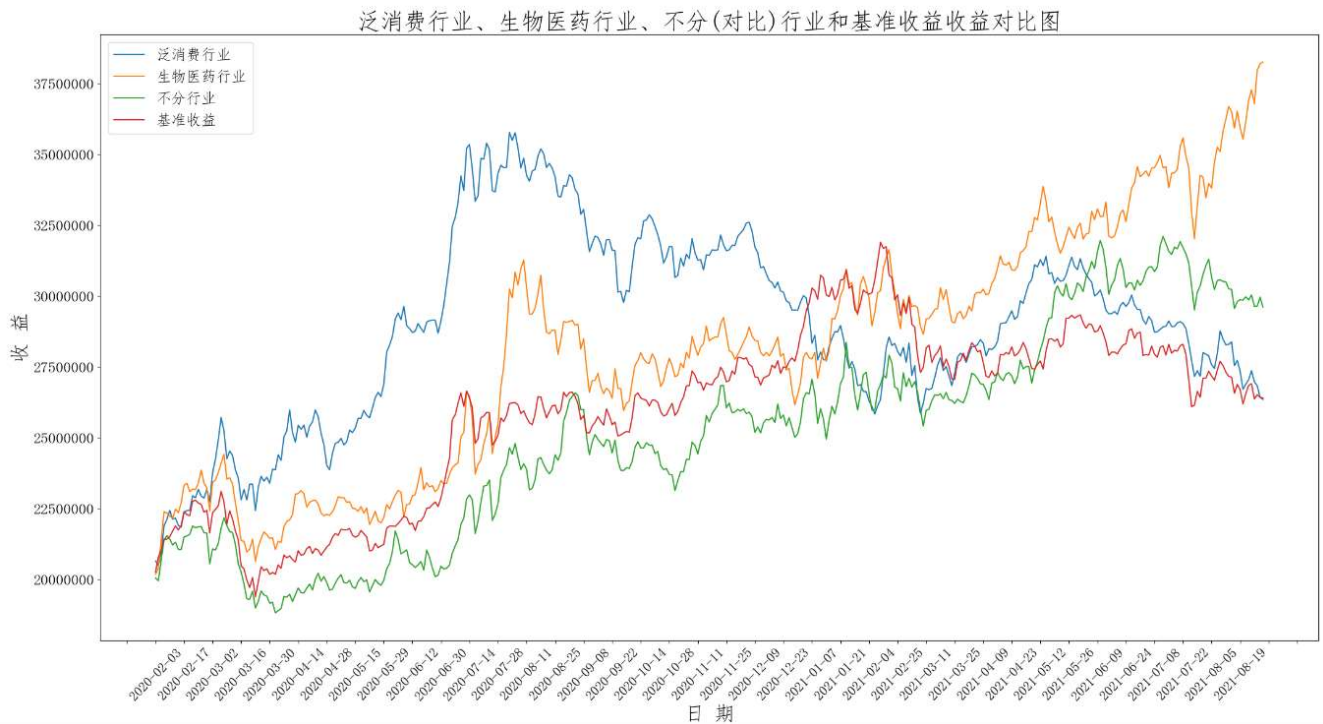


图 3-9 所有结果收益对比图

表 3-5 评价指标

	收益指标	
	年化收益	最大回撤
泛消费行业	17.63%	27.77%
生物医药行业	53.17%	17.00%
不分行业（对比）	29.64%	15.13%
基准收益	19.29%	18.19%

由图 3-9 可知，泛消费行业和生物医药行业在回测区间上，整体的走势都要比不分行业的好。同时我们将泛消费行业、生物医药行业和不分行业的收益分别与基准收益进行对比，具体可以分为走势对比和收益指标对比两个方面，其中收益指标本次只选取年化收益率和最大回测两个指标，因为选股评估中回测阶段并没有使用任何择时策略，两种对比的具体结果为下文所述：

1. 走势对比, 我们可以看到泛消费行业整体优于基准收益, 经计算泛消费行业大约有 84.06% 的时间收益高于基准收益, 生物医药行业则是有 87.15% 的时间的收益高于基准收益, 而不分行业则仅有 22.11% 的时间收益高于基准收益, 即不分行业的收益大部分时间低于基准收益。

2. 收益指标对比, 从表 3-5 中可以看到生物医药行业的年化收益率比基准收益和不分行业的都要高, 同时生物医药行业的最大回撤却只比不分行业的大, 但比基准收益的要小。而泛消费行业的年化收益率比基准收益和不分行业的都要小, 泛消费行业的最大回撤却比基准和不分行业的都要大。

综合上面两方面的对比, 我们可以知道, 生物医药行业所选取的股票无论是走势还是收益指标都远远的优于基准收益和不分行业。而对于泛消费行业来说, 如果只单单看指标, 那么泛消费行业选取的股票的收益比基准收益和不分行业的收益都要差。但是根据走势对比的结果可以发现泛消费行业的收益大部分时间都是高于基准收益和不分行业, 这说明泛消费行业实际上收益还是比较好。通过泛消费行业的指标和走势得出截然相反的结果可知, 交易时机的不合理有可能会导出选取了走势较好的股票, 却无法获得较好收益的情况出现, 所以合理的交易策略对最终的收益具有很大的影响。

第四节 构建择时策略与回测验证

一、计算长期均线

本文分别对泛消费行业和生物医药行业选出的每只股票选取出最佳的长期均线, 并且以滚动的方式对每只股票的长期均线进行更新。其中具体的操作为下文所述:

1. 对每只月股票分配 20 万资金, 以择时月份第一个交易日往前推七年的时间作为获取长期均线的训练时间, 以收盘价上穿或者下穿长期均线作为买入和卖出信号, 对股票进行买入和卖出的交易。其中具体信号为如下描述:

a. 当日收盘价大于当日均线值, 昨日收盘价小于昨日均线值, 发出买入信号。

b. 当日收盘价小于当日均线值, 昨日收盘价大于昨日均线值, 发出卖出信号。

2. 按照上述的方式将均线值从 120 到 240 逐个取值, 计算每个均线值获取的收益, 选取收益最高所对应的均线值作为股票当月的最佳长期均线值。

3. 通过滚动的方式对股票的最佳长期均线进行更新, 即在下月的第一个交易日往前推七年作为新的长期均线训练时间。

二、训练深度强化学习 DQN 算法

（一）模型的训练

对于 DQN 算法的训练，本文综合考虑 DQN 算法主要用于短期择时以及股票市场环境的复杂性以及行业特点等三个因素，对不同的行业采取不同长度的训练数据对模型进行训练。

泛消费行业的股票具有流通快、波动大等特点。同时泛消费行业因为涉及的面比较广，以至于泛消费行业的股票市场环境复杂性比较高，所以本文将选取半年的训练数据对模型进行训练。即如果 2020 年 2 月是当前择时月，则使用 2019 年 8 月到 2020 年 1 月的数据进行训练。同时采用滚动的方式进行训练，以月为单位，每个自然月都重新使用前 6 个月的数据训练 DQN 算法。

生物医药行业的股票则受到医药行业研发周期长，技术含量高等特点的影响，其股票价格波动幅度并不剧烈。而生物医药行业则涉及的范围相对于泛消费行业更小，所以生物医药行业的股票市场复杂性并不高，所以本文将选取一年的训练数据对模型进行训练。即如果 2020 年 2 月是当前择时月，则使用 2019 年 2 月到 2020 年 1 月的数据进行训练。同时采用滚动的方式进行训练，以月为单位，每个自然月都重新使用前一年的数据训练 DQN 算法。

具体的模型执行流程如图 3-10 所示：

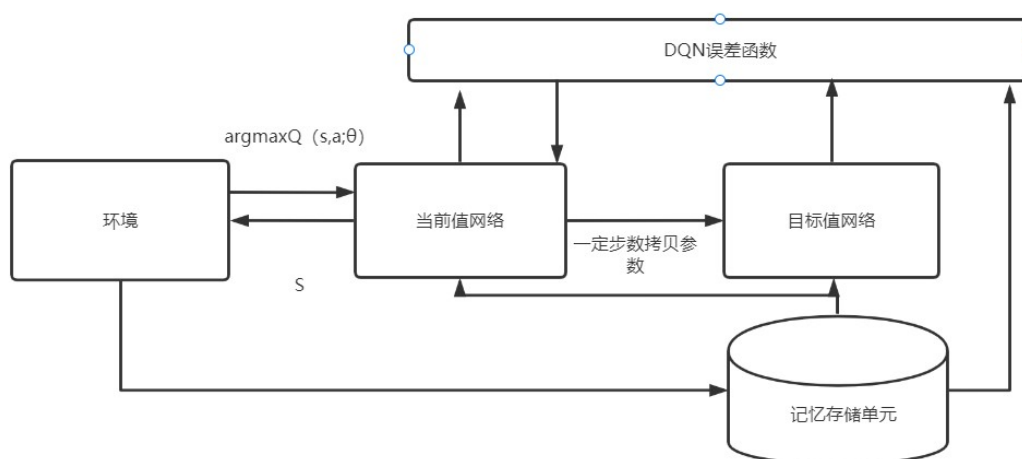


图 3-10 DQN 算法训练流程

三、择时策略实证结果评估

将泛消费行业和生物医药行业经过多因子选股模型选取的股票在 2020 年 2 月 3 日至 2021 年 8 月 31 日上结合择时策略进行回测，并假定股票的交易是零滑点¹⁴。而进行回测时具体的调仓策略详见第二章。对泛消费行业和生物医药行业选取的股票

¹⁴ 零滑点：预期买入、卖出的价格和实际买入、卖出的价格一致，不存在偏差称为零滑点。

在回测区间上使用新型择时交易策略进行交易，并将择时交易的结果和为择时交易的结果进行对比分析，其中择时和未择时收益的对比如图 3-11 和 3-12 所示：

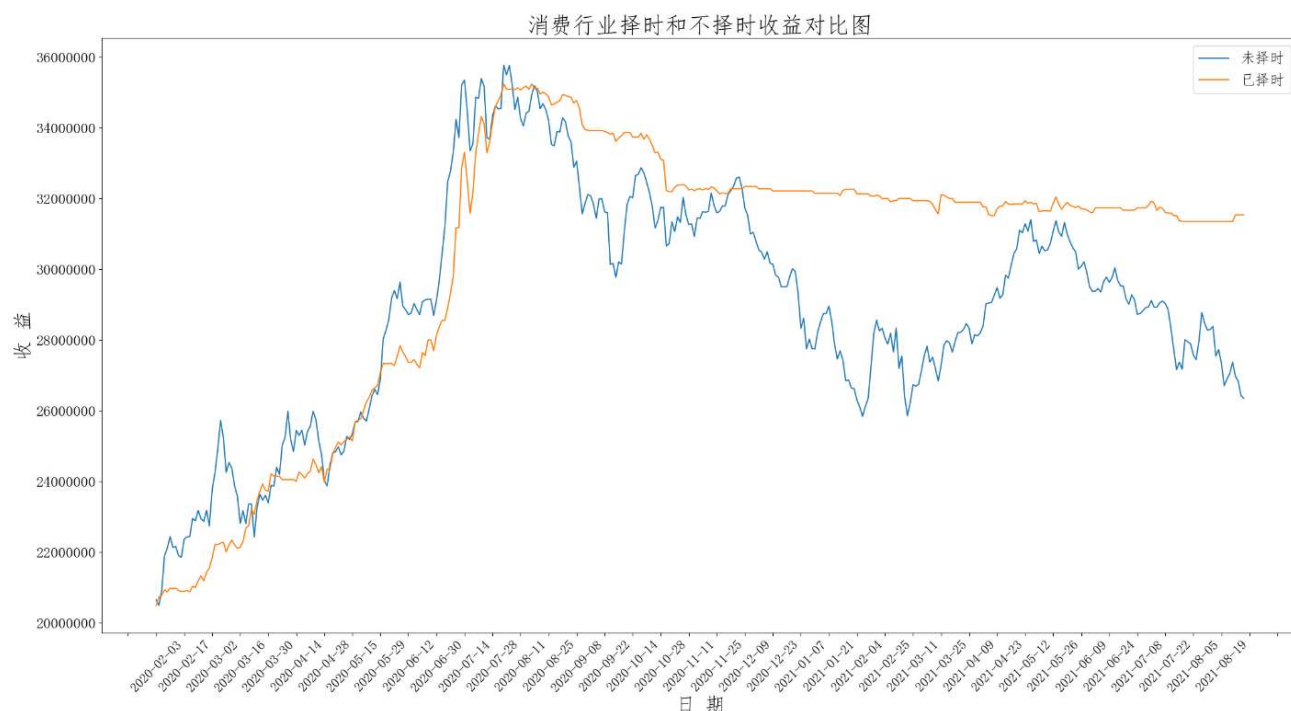


图 3-11 泛消费行业择时和未择时收益对比图



图 3-12 生物医药行业择时未择时收益对比图

表 3-6 评价指标

	收益指标			
	年化收益	超额年化收益	最大回撤	夏普比率
泛消费行业（已择时）	33.35%	11.79%	11.04%	0.56
泛消费行业（未择时）	17.63%	-1.41%	27.77%	-0.05
生物医药行业（已择时）	59.51%	32.86%	11.03%	1.51
生物医药行业（未择时）	53.17%	28.34%	17.00%	1.35

从图 3-11 至 3-12 以及表 3-6 的收益对比可以发现，择时策略带来的收益比较好，尤其当股票走势下降的时候，择时策略比较敏锐的捕捉到了，从而缩小了收益的回撤¹⁵，提升了整体收益。从整体上看使用了择时策略的分组获取的收益一直保持在较高的水平，没有出现收益为负的状态，收益较为稳定。

根据表 3-6 中泛消费行业的各项指标可以发现，使用择时策略的年化收益几乎是没使用择时策略的两倍，于此同时使用择时策略的最大回撤却只有未使用择时策略的 1/3 左右。这很好的证明了择时策略成功的预判了股票的走势，做出了正确的交易。而对于生物医药行业来说，虽然使用择时策略和没有使用择时策略带来的提升并不明显，但是通过使用择时策略在一定程度上提升了收益。最后通过两个行业的择时和不择时超额收益率进行对比可以明显看到，择时策略有效的提升了 Alpha 收益¹⁶，而对比夏普比率则可以明显的发现策略提升了夏普比率，尤其是泛消费行业把未择时情况下负的夏普比率提升到了正的夏普比率，有效的提升了相等风险下的收益。

综上所述，长期均线和深度强化学习结合的新型择时策略有效的提升了收益，所以这种择时策略是有效的。

¹⁵ 回撤：指的是某一段时间内股票的净值从最高点开始回落到最低点的幅度。

¹⁶ Alpha 收益：个股收益与平均收益之间的差值称为 Alpha 收益，即高于平均收益的部分。

第四章 总结与展望

第一节 总结

本文的主要的研究有两个部分：第一个部分是基于行业构建多因子选股模型；第二个部分是对分行业构建的多因子选股模型选取的股票设计一种择时交易策略的研究。

第一个部分，基于行业构建多因子选股模型主要做了如下几个方面的工作：

1.通过聚宽量化数据平台获取了泛消费行业和生物医药行业共计 481 只股票于 2012 年 1 月 4 日到 2021 年 8 月 31 日期间的包含开盘价、收盘价、最高价、最低价和成交量等五个维度的数据。通过通联量化交易平台，获取了七个大类共计 232 个因子在 2012 年 1 月 4 日到 2021 年 8 月 31 日期间的日度数据，共计 1.6GB。

2.对股票数据进行处理，首先对股票进行异常检查，一共剔除了 30 只不符合要求的股票。其次对股票数据以月为单位，以每日收盘价计算股票月内的总涨幅和所属行业同一个月份的月内涨幅进行比较生成数据的标签，通过比较一共生成 51865 个标签值。然后对股票数据和因子数据进行标准化处理。最后对因子进行选取，通过信息系数对因子进行初步筛选，将因子和未来收益之间的关系不显著的因子给剔除掉。紧接着使用自己构建的最终因子筛选方法对因子进行刷选，最后一共对泛消费行业、生物医药行业和不分行业三个行业选取了 12、16 和 12 个因子用于构建多因子选股模型。

分行业构建多因子选股模型，以行业为单位将选取的因子作为输入的特征输入到 LightGBM、XGBoost 和 AdaBoost 等三个算法中构建多因子选股模型。并且以月度为单位进行滚动训练，其中训练数据集的长度为 96 个月数据，测试周期数据共有 20 个月数据。

根据上述步骤完成实验之后，通过最终得出的结果可以发现，因“股”制宜，基于行业的多因子选股模型可以准确的选取出未来一个月内走势好的股票。

第二个部分，构建交易策略。本文选取的是长期均线和强化学习融合构建新型择时交易策略，主要做了如下几个方面的工作：

1.对基于行业的构建的多因子模型选取的所有股票都求取上七个年度中的最佳长期均线，其中长期均线的取值范围选取的 120 到 240 日。对于每个行业的 19 期股票以月为单位滚动求长期均线，每个行业共计求得 190 个最佳长期均线值。

2.对股票的开盘价、收盘价、最高价、最低价和成交量等数据标准化，然后对股票市场实现 DQN 算法，并结合行业特点为泛消费行业选取了半年、生物医药行业选取了一年的数据进行训练，其中每个行业的 19 期股票共计训练得到 190 个模型。

3.构建使用长期均线 and 强化学习融合信号的回测交易系统，并对每个行业选取的股票于 2020 年 1 月 2 日到 2021 年 8 月 31 日上使用历史数据进行回测。

使用上述择时策略，完成回测之后，根据得出的结果可以发现，择时交易策略可以敏锐的捕捉到股票走势的变化，及时的调整交易行为，降低了收益回撤的幅度，稳定了最终的股票收益。

第二节 不足与展望

本文基于行业构建了多因子选股模型和择时策略，虽然从最后的结果来看基于行业的多因子选股和择时策略是一种有效提升收益的方法，但是本文的研究仍存在不足之处。比如由于受到实验时间有限、论文写作时间有限和个人自身等原因导致本文实证的案例并不是很丰富，所以未来的研究中可以补充案例用于进一步论证本文的研究结果。同时本文研究在未来还可以进一步进行拓展，主要是分为股票类别的划分、候选因子、算法模型和择时策略等四个方面的拓展。

首先，股票类别的划分。对于股票类别的划分在未来可以通过使用其他的分类方式对股票进行分类，比如使用行业生命周期的初创阶段、成长阶段、成熟阶段和衰退阶段等四个阶段对股票进行分类，并以此来构建基于股票类别的多因子选股模型。

其次，候选因子。对于候选因子，本文使用的是通联量化数据平台上的因子，未来的研究中可以自己股票市场对股票进行数据挖掘，从海量的数据中挖掘出市场中蕴含的更加深层次的关系，更进一步揭示影响股票未来一段时间内走势的因素，从而进一步提高选股模型选取走势较好的股票的概率。

然后，算法模型。对于构建模型的算法本文选取的是 LightGBM、XGBoost 和 AdaBoost 等三种算法，未来构建多因子选股模型的时候可以尝试其他的方法，比如使用深度学习来构建多因子选股模型，通过神经网络来刻画当前因子值和股票未来收益之间的关系。

最后，择时策略。对于择时策略来说，本文选取的是长期均线 and 强化学习结合构建择时策略，在未来的研究中可以通过更换强化学习的算法，使用更加先进的强化学习算法来构建这种复合信号。同理，我们也可以使用其他的传统指标与强化学习进行结合，从而产生其他的择时交易策略，用于股票市场的择时交易。

参考文献

- [1]王宁.基于创业板行业分层的机器学习选股策略研究[D].上海师范大学,2021.
- [2]宋安娜. 财务柔性、技术创新与企业价值[D].浙江大学,2021.
- [3]李雨诗. 基于集成学习与深度学习的文本情感分析研究[D].兰州大学,2021.
- [4]林永峰,沈彦,李禹汉,陈桦.基于择时的多因子选股模型[J].信息技术与标准化,2021(06):44-50.
- [5]袁晨光. 基于投票集成学习算法的多因子量化选股方案研究[D].上海师范大学,2021.
- [6]马瑞雪.组合滤波与希尔伯特变换的短线交易择时研究[J].时代金融,2021(13):76-80.
- [7]黄蕊. 基于 LightGBM 算法的行业轮动多因子选股方案研究[D].上海师范大学,2021.
- [8]高子宜. 基于随机森林的股价走势预测研究[D].中国政法大学,2021.
- [9]石倩.基于集成学习的股价涨跌预测与交易策略研究[D].中南财经政法大学,2021.
- [10]董焕彬. 基于多智能体强化学习模型在 A 股择时和选股应用研究[D].浙江大学,2020.
- [11]于龙飞. 基于深度学习的股市量化交易系统设计与实现[D].山东大学,2020.
- [12]葛槽漠,周显.基于 XGBoost 的多因子选股模型[J].信息技术与标准化,2020(05):36-41.
- [13]汪鲁越. 生物医药行业投资价值研究[D].浙江大学,2020.
- [14]常太星. 基于多因子选股和隐马尔科夫模型择时的量化策略研究[D].哈尔滨工业大学,2020.
- [15]何路. 多因子量化选股及投资者情绪择时策略的实证检验[D].南京大学,2020.
- [16]祝养豹. 基于 XGBoost 和 LightGBM 算法的多因子选股方案设计[D].南京大学,2020.
- [17]欧阳明哲.基于 GRU 神经网络的多因子量化选取策略[D].中南财经政法大学,2021.
- [18]李一帆. 我国生物医药行业上市公司投资价值分析[D].云南财经大学,2020.
- [19]王一卓. 基于 Boosting 算法的多因子量化选股实证研究[D].山东大学,2020.
- [20]腾讯证券.2020 年中国股民行为报告[R].
- [21]廖安东. 基于集成学习算法的 A 股投资策略研究[D].电子科技大学,2020.
- [22]李奉珂. 基于投资者情绪的多因子选股模型实证研究[D].西南财经大学,2019.
- [23]韩立强. 基于 boosting 模型的逆向投资策略研究[D].浙江大学,2019.
- [24]华宇. 基于 BP 神经网络模型的股票择时研究[D].沈阳工业大学,2019.
- [25]李杨. 基于 XGBoost 的基本面量化模型[D].山东大学,2019.
- [26]霍丽佳. 基于 AdaBoost 算法多因子选股模型的应用研究[D].华中科技大学,2019.

- [27]姜加才. 基于 LightGBM 算法的量化选股策略方案策划[D].上海师范大学,2019.
- [28]张冬阳. 基于 Logistic 回归的 Barra 因子选股模型研究[D].南京大学,2018.
- [29]姚瞳彤. AdaBoost 算法在量化投资中的改进与应用研究[D].暨南大学,2018.
- [30]王凯. 基于集成学习的量化选股策略研究[D].华南理工大学,2017.
- [31]李想. 基于 XGBoost 算法的多因子量化选股方案策划[D].上海师范大学,2017.
- [32]刘全,翟建伟,章宗长,钟珊,周倩,章鹏,徐进.深度强化学习综述[J].计算机学报,2018,41(01):1-27.
- [33]夏烈阳. 大数据背景下基于 Web 日志的用户访问模式挖掘研究[D].云南财经大学,2019.
- [34]张子薇. 基于改进 GBDT 算法的光伏发电功率预测研究[D].华北电力大学,2018.
- [35]Liu Meng,Luo Kaiping,Zhang Junhuan,Chen Shengli. A stock selection algorithm hybridizing grey wolf optimizer and support vector regression[J]. Expert Systems With Applications,2021,179:
- [36]Koshiyama Adriano,Firoozye Nick,Treleaven Philip. Generative adversarial networks for financial trading strategies fine-tuning and combination[J]. Quantitative Finance,2021,21(5):
- [37]Yihua Zhong , Lan Luo , Xinyi Wang , Jinlian Yang. Multi-factor Stock Selection Model Based on Machine Learning[J]. Engineering Letters,2021,29(1):
- [38]Budiharto Widodo. Data science approach to stock prices forecasting in Indonesia during Covid-19 using Long Short-Term Memory (LSTM)[J]. Journal of big data,2021,8(1):
- [39]Nguyen Nguyet,Nguyen Dung. Global Stock Selection with Hidden Markov Model[J]. Risks,2020,9(1):
- [40]Huotari Tommi,Savolainen Jyrki,Collan Mikael. Deep Reinforcement Learning Agent for S&P 500 Stock Selection[J]. Axioms,2020,9(4):
- [41]Wu Xing,Chen Haolei,Wang Jianjia,Troiano Luigi,Loia Vincenzo,Fujita Hamido. Adaptive stock trading strategies with deep reinforcement learning methods[J]. Information Sciences,2020,538:
- [42]Hyungjun Park,Min Kyu Sim,Dong Gu Choi. An intelligent financial portfolio trading strategy using deep Q-learning[J]. Expert Systems With Applications,2020,158:
- [43]Iwao Maeda,David deGraw,Michiharu Kitano,Hiroyasu Matsushima,Hiroki Sakaji,Kiyoshi Izumi,Atsuo Kato. Deep Reinforcement Learning in Agent Based Financial Market Simulation[J]. Journal of Risk and Financial Management,2020,13(4):
- [44]Zheng Tan,Ziqin Yan,Guangwei Zhu. Stock selection with random forest: An exploitation of excess return in the Chinese stock market[J]. Heliyon,2019,5(8):
- [45]Zhige Li, Derek Yang, Li Zhao, Jiang Bian, Tao Qin, and Tie-Yan Liu. 2019. Individualized Indicator for All: Stock-wise Technical Indicator Optimization with

- Stock Embedding. In The 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '19), August 4–8, 2019, Anchorage, AK, USA. ACM, New York, NY, USA, 9 pages.
- [46]Keywan Christian Rasekhschaffe,Robert C. Jones. Machine Learning for Stock Selection[J]. Financial Analysts Journal,2019,75(3):
- [47]Elena Krashennnikova,Javier García,Roberto Maestre,Fernando Fernández. Reinforcement learning for pricing strategy optimization in the insurance industry[J]. Engineering Applications of Artificial Intelligence,2019,80:
- [48]Masanori Hirano,Hiroki Sakaji,Shoko Kimura,Kiyoshi Izumi,Hiroyasu Matsushima,Shintaro Nagao,Atsuo Kato. Related Stocks Selection with Data Collaboration Using Text Mining †[J]. Information,2019,10(3):
- [49]Multi-factor Stock Selection Model Based on Adaboost[J]. Business and Economic Research,2018,8(4):
- [50]Deng Yue,Bao Feng,Kong Youyong,Ren Zhiquan,Dai Qionghai. Deep Direct Reinforcement Learning for Financial Signal Representation and Trading[J]. IEEE transactions on neural networks and learning systems,2017,28(3):

附录 1

证监会泛消费行业分类代码说明

门类	大类	类别名称	说明
C	13	农副食品加工业	指直接以农、林、牧、渔业产品为原料进行的谷物磨制、饲料加工、植物油和制糖加工、屠宰及肉类加工、水产品加工，以及蔬菜、水果和坚果等食品的加工
	14	食品制造业	
	15	酒、饮料和精制茶制造业	
	16	烟草制品业	
	17	纺织业	
	18	纺织服装、服饰业	
	19	皮革、毛皮、羽毛及其制品和制鞋业	
	20	木材加工和木、竹、藤、棕、草制品业	
	21	家具制造业	指用木材、金属、塑料、竹、藤等材料制作的，具有坐卧、凭倚、储藏、间隔等功能，可用于住宅、旅馆、办公室、学校、餐馆、医院、剧场、公园、船舰、飞机、机动车等任何场所的各种家具的制造
	22	造纸和纸制品业	
	23	印刷和记录媒介复制业	
	24	文教、工美、体育和娱乐用品制造业	
F	51	批发业	指向其他批发或零售单位（含个体经营者）及其他企事业单位、机关团体等批量销售生活用品、生产资料的活动，以及从事进出口贸易和贸易经纪与代理的活动，包括拥有货物所有权，并以本单位（公司）的名义进行交易活动，也包括不拥有货物的所有权，收取佣金的商品代理、商品代售活动；本类还包括各类商品批发市场中固定摊位的批发活动，以及以销售为目的的收购活动
	52	零售业	指百货商店、超级市场、专门零售商店、品牌专卖店、售货摊等主要面向最终消费者（如居民等）的销售活动，以互联网、邮政、电话、售货机等方式的销售活动，还包括在同一地点，后面加工生产，前面销售的店铺（如面包房）；谷物、种子、饲料、牲畜、矿产品、生产用原料、化工原料、农用化工产品、机械设备（乘用车、计算机及通信设备除外）等生产资料的销售不作为零售活动；多数零售商对

			其销售的货物拥有所有权，但有些则是充当委托人的代理人，进行委托销售或以收取佣金的方式进行销售
H	61	住宿业	指为旅行者提供短期留宿场所的活动，有些单位只提供住宿，也有些单位提供住宿、饮食、商务、娱乐一体的服务，本类不包括主要按月或按年长期出租房屋住所的活动
	62	餐饮业	指通过即时制作加工、商业销售和服务性劳动等，向消费者提供食品 and 消费场所及设施的服务
P	82	教育	

证监会生物医药行业分类代码说明

门类	大类	类别名称	说明
C	26	化学原料和化学制品制造业	
	27	医药制造业	
Q	83	卫生	

附录 2

泛消费行业初始股票池

特力 A	飞亚达	国药一致	深纺织 A	海王生物	广聚能源	中成股份	常山北明
国际实业	英特集团	民生控股	合肥百货	通程控股	华天酒店	晨鸣纸业	鄂武商 A
ST 国医	岭南控股	广弘控股	泰山石油	我爱我家	泸州老窖	太阳能	德龙汇能
古井贡酒	远大控股	三木集团	西王食品	泰达股份	厦门信达	正虹科技	浙江震元
中兴商业	西安饮食	鲁泰 A	燕京啤酒	漳州发展	中百集团	居然之家	凯撒旅业
创维数字	陕西金叶	德展健康	天音控股	粤桂股份	承德露露	华茂股份	高鸿股份
五粮液	顺鑫农业	张裕 A	新希望	中嘉博创	双汇发展	浙商中拓	大亚圣象
兰州黄河	华东医药	银泰黄金	越秀金控	中国中期	伟星股份	凯恩股份	苏宁易购
七匹狼	旺能环境	联创电子	华孚时尚	兔宝宝	景兴纸业	太阳纸业	孚日股份
新野纺织	江苏国泰	天康生物	广博股份	东港股份	天邦股份	宏达高科	报喜鸟
正邦科技	游族网络	全聚德	广百股份	劲嘉股份	如意集团	三全食品	东华能源
合兴包装	鸿博股份	安妮股份	奥特佳	盛新锂能	步步高	恩华药业	新华都
美邦服饰	友阿股份	保龄宝	星期六	奥飞娱乐	罗莱生活	齐心集团	美盈森
洋河股份	海大集团	富安娜	皇氏集团	得利斯	高乐股份	大北农	维信诺
联发股份	梦洁股份	嘉欣丝绸	爱施德	天虹股份	凯撒文化	众业达	珠江啤酒
嘉事堂	双塔食品	雪松发展	嘉麟杰	浙江永强	华斯股份	佳隆股份	搜于特
涪陵榨菜	中顺洁柔	金字火腿	旷达科技	恺英网络	齐峰新材	探路者	吉峰科技
星辉娱乐	万顺新材	恒信东方	英唐智控	晨光生物	汤臣倍健	睿智医药	宁波联合
中国医药	五矿发展	古越龙山	国投资本	浙江富润	上海梅林	湘财股份	青山纸业
美尔雅	浙江东日	浙江东方	金健米业	弘业股份	重庆啤酒	浪莎股份	维科技术
建发股份	华升股份	雅戈尔	泉阳泉	伊力特	金种子酒	江苏阳光	金鹰股份
圆通速递	民丰特纸	冠农股份	首旅酒店	开开实业	嘉化能源	东方创业	江苏舜天
广汇汽车	安琪酵母	恒顺醋业	华泰股份	大东方	美克家居	恒丰纸业	华联综超
三房巷	北巴传媒	海澜之家	红豆股份	三元股份	冠豪高新	通威股份	瑞贝卡
华纺股份	福能股份	凤竹纺织	中化国际	国药股份	贵州茅台	莫高股份	老白干酒
山鹰国际	惠泉啤酒	光明乳业	青岛啤酒	汇通能源	金枫酒业	申达股份	新世界
龙头股份	外高桥	申华控股	豫园股份	昂立教育	南京新百	东百集团	大商股份
欧亚集团	物产中大	南宁百货	南京医药	首商股份	重庆百货	中国高科	新华锦
中粮糖业	丽尚国潮	辽宁成大	东方银星	锦江酒店	厦门国贸	汉商集团	新华百货
山西汾酒	东方集团	杭州解百	上海物贸	益民集团	兰生股份	百联股份	茂业商业
香溢融通	第一医药	上海九百	海欣股份	宁波中百	银座股份	王府井	北京城乡
百大集团	中炬高新	梅花生物	伊利股份	岳阳林纸	博汇纸业	汇鸿集团	航民股份
马应龙	九州通	金陵饭店	上海医药	际华集团	永辉超市	*ST 大集	大连友谊
天娱数科	鼎龙文化	*ST 济堂	ST 宏图	*ST 华资	ST 时万	中央商场	ST 维维
*ST 广珠	ST 海越	*ST 绿庭	舍得酒业	金开新能	上海易连	深粮控股	

生物医药行业初始股票池

丽珠集团	四环生物	云南白药	安道麦 A	渝三峡 A	海南海药	东北制药	吉林敖东
天茂集团	英力特	仁和药业	长春高新	远兴能源	四川美丰	普洛药业	新华制药
通化金马	北大医药	甘肃电投	鲁西化工	华特达因	金陵药业	中粮科技	广济药业
大庆华科	九芝堂	诚志股份	华润三九	新和成	华邦健康	华兰生物	传化智联
京新药业	科华生物	达安基因	保利联合	双鹭药业	云南能投	德美化工	中钢天源
黑猫股份	中泰化学	青岛金王	南岭民爆	海翔药业	沃华医药	紫鑫药业	湘潭电化
安纳达	红宝丽	莱茵生物	东方锆业	芭田股份	嘉应制药	宏达新材	诺普信
江南化工	北化股份	联化科技	上海莱士	利尔化学	华昌化工	桂林三金	奇正藏药
信立泰	众生药业	永太科技	仙琚制药	新纶新材	精华制药	同德化工	神剑股份
永安药业	亚太药业	天原股份	信邦制药	长青股份	力生制药	海普瑞	多氟多
齐翔腾达	雅克科技	延安必康	汉森制药	科伦药业	贵州百灵	太安堂	誉衡药业
闰土股份	龙星化工	华软科技	百川股份	天齐锂业	宝莫股份	雅化集团	莱美药业
安科生物	北陆药业	硅宝科技	红日药业	新宙邦	上海凯宝	回天新材	福瑞股份
鼎龙股份	天龙集团	安诺其	三聚环保	奥克股份	康芝药业	建新股份	吉药控股
新开源	华仁药业	瑞普生物	阳谷华泰	智飞生物	青松股份	宝利国际	沃森生物
香雪制药	华润双鹤	人福医药	同仁堂	云天化	太极集团	乐凯胶片	兴发集团
巨化股份	天坛生物	中牧股份	复星医药	生物股份	西藏药业	浙江医药	太龙药业
圣济堂	两面针	中恒集团	海正药业	恒瑞医药	亿利洁能	羚锐制药	万华化学
上海家化	中盐化工	中新药业	白云山	亚宝药业	浙江龙盛	红星发展	昊华科技
健康元	江山股份	三友化工	国药现代	昆药集团	华鲁恒升	片仔癀	蓝光发展
千金药业	双良节能	扬农化工	天药股份	联环药业	华海药业	天士力	康缘药业
康恩贝	益佰制药	新安股份	神奇制药	氯碱化工	华谊集团	哈药股份	湖南海利
江苏索普	江中药业	鲁抗医药	钱江生化	新奥股份	华北制药	人民同泰	通化东宝
健民集团	淮北矿业	滨化股份	*ST 浪奇	ST 红太阳	*ST 恒康	兆新股份	乐通股份
ST 金正	*ST 澄星	*ST 环球	ST 瀚叶	六国化工	交大昂立	ST 中珠	ST 榕泰
ST 熊猫	东阿阿胶						

附录 3

初始候选因子			
ILLIQUIDITY	GrossIncomeRatio	MA10RegressCoeff12	BIAS5
收益相对金额比	销售毛利率	10 日价格平均线 12 日 线性回归系数	5 日乖离率
IntangibleAssetRatio	DDNBT	ROE	ChaikinVolatility
无形资产比率	下跌贝塔	权益回报率	佳庆离散指标
EGRO	plusDI	ARTDays	PE
5 年收益增长率	上升指标	应收账款周转天数	市盈率
ChaikinOscillator	Aroon	MA10RegressCoeff6	AccountsPayablesTRate
佳庆指标	趋势反转的变化	10 日价格平均线 6 日线 性回归系数	应付账款周转率
minusDI	BLEV	CTP5	NPTtoTOR
下降指标	账面杠杆	5 年平均现金流市值比	净利润与营业总收入 之比
ARC	CFO2EV	VOL10	DebtEquityRatio
变化率指数均值	经营活动产生的现金流 流量净额与企业价值 之比	10 日平均换手率	产权比率
InventoryTDays	VEMA12	EMA120	DVRAT
存货周转天数	成交量的 12 日指数移 动平均	120 日指数移动均线	收益相对波动
CCI5	EMV14	RC12	ADXR
5 日顺势指标	简易波动指标	12 日变化率指数	相对平均动向指数
ROC20	Skewness	TEMA10	ROA5
20 日变动速率	股价偏度	10 日三重指数移动平均 线	5 年资产回报率
VOL60	DDNSR	GREV	REVS20
60 日平均换手率	下跌波动	分析师盈利预测变化趋 势	股票的 20 日收益
OBV20	OperatingRevenueGro wRate	ATR6	FinancingCashGrowRate
20 日能量潮指标	营业收入增长率	6 日均幅指标	筹资活动产生的现金流 流量净额增长率
VSTD20	MoneyFlow20	FiftyTwoWeekHigh	KDJ_K
20 日成交量标准 差	20 日资金流量	当前价格处于过去 1 年 股价的位置	随机指标
PVI	MA10Close	KDJ_D	OperatingExpenseRate
正成交量指标	均线价格比	随机指标	营业费用与营业总收入 之比
ARBR	TotalProfitCostRatio	AR	FinancialExpenseRate
人气指标	成本费用利润率	人气指标	财务费用与营业总收入 之比
BackwardADJ	VR	RSTR12	OperatingProfitToTOR
股价向后复权因 子	成交量比率	12 月相对强势	营业利润与营业总收入 之比

BIAS20 20 日乖离率	CashToCurrentLiability 现金比率	EMA5 5 日指数移动均线	TVSTD20 20 日成交金额的标准差
REVS10 股票的 10 日收益	PSY 心理线指标	OperCashInToCurrentLiability 现金流流动负债比	InvestCashGrowRate 投资活动产生的现金流量净额增长率
SUOI 未预期毛利	PS 市销率	CTOP 现金流市值比	AdminiExpenseRate 管理费用与营业总收入之比
FY12P 分析师盈利预测	BR 意愿指标	VOL20 20 日平均换手率	BBI 多空指数
Volatility 换手率相对波动率	TotalAssetsTRate 总资产周转率	TRIX10 10 日收盘价三重指数平滑移动平均指标	NOCFTToOperatingNI 经营活动产生的现金流量净额与经营活动净收益之比
JDQS20 阶段强势指标	SRMI 修正动量指标	GSREV 分析师盈收预测变化趋势	OBV6 6 日能量潮指标
VROC12 12 日量变动速率指标	ETOP 收益市值比	EMA10 10 日指数移动均线	VROC6 6 日量变动速率指标
UOS 终极指标	TEMA5 5 日三重指数移动平均线	ACCA 现金流资产比和资产回报率之差	MA60 60 日移动均线
AD 累积/派发线	SaleServiceCashToOR 销售商品提供劳务收到的现金与营业收入之比	FixAssetRatio 固定资产比率	EMV6 简易波动指标
DebtsAssetRatio 债务总资产比	RSI 相对强弱指标	SalesCostRatio 销售成本率	PB 市净率
DHILO 波幅中位数	Hurst 赫斯特指数	BondsPayableToAsset 应付债券与总资产之比	CurrentAssetsRatio 流动资产比率
MLEV 市场杠杆	FSALESG 未来预期盈收增长	EPS 基本每股收益	OBV 能量潮指标
TaxRatio 销售税金率	SUE 未预期盈余	TRIX5 5 日收盘价三重指数平滑移动平均指标	EMA60 60 日指数移动均线
DAVOL5 5 日平均换手率与 120 日平均换手率	DilutedEPS 稀释每股收益	PVT6 因子 PVT 的 6 日均值	
HSIGMA 历史波动	MFI 资金流量指标	KlingerOscillator 成交量摆动指标	PLRC12 12 日收盘价格价格线性回归系数
EBITToTOR 息税前利润与营业总收入之比	CCI20 20 日顺势指标	Ulcer5 5 日向下的波动性	OperatingProfitRatio 营业利润率
DDNCR 下跌相关系数	MA10 10 日移动均线	BullPower 多头力道	CurrentRatio 流动比率
GREC	MTM	ROA	LFLO

分析师推荐评级 变化趋势	动量指标	资产回报率	对数流通市值
CoppockCurve	VEMA5	MAWVAD	DBCD
估波指标	成交量的 5 日指数移 动平均	因子 WVAD 的 6 日均 值	异同离差乖离率
ROE5	VEMA26	REVS5	DDI
5 年权益回报率	成交量的 26 日指数移 动平均	股票的 5 日收益	方向标准离差指数
NetAssetGrowRate	EARNMOM	LCAP	CCI10
净资产增长率	八季度净利润变化趋 势	对数市值	10 日顺势指标
DAVOL10	DAREC	NPParentCompanyGrow Rate	MassIndex
10 日平均换手率 与 120 日平均换 手率	分析师推荐评级变化	归属母公司股东的净利 润增长率	梅斯线
BBIC	VOL240	ASI	TotalAssetGrowRate
因子 BBI 除以收 盘价得到	240 日平均换手率	累计振动升降指标	总资产增长率
DAREV	REC	PVT12	MA20
分析师盈利预测 变化	分析师推荐评级	因子 PVT 的 12 日均值	20 日移动均线
ACD6	DAVOL20	ROC6	CR20
6 日收集派发指标	20 日平均换手率与 120 日平均换手率	6 日变动速率	CR 指标以上一个计算 周期
PCF	DEGM	Elder	NetProfitRatio
市现率	毛利率增长	艾达透视指标	销售净利率
VOL5	BIAS60	RC24	CMRA
5 日平均换手率	60 日乖离率	24 日变化率指数	24 月累计收益
FEARNG	EMA12	ARTRate	MACD
未来预期盈利增 长	12 日指数移动均线	应收账款周转率	平滑异同移动平均线
VMACD	NVI	BollDown	TA2EV
成交量量指数平 滑异同移动平均 线	负成交量指标	下轨线	资产总计与企业价值 之比
CurrentAssetsTRate	AccountsPayablesTDay s	PVT	AD6
流动资产周转率	应付账款周转天数	价量趋势	累积/派发线
SwingIndex	WVAD	ACD20	BearPower
振动升降指标	威廉变异离散量	20 日收集派发指标	空头力道
OperCashGrowRate	MTMMA	EMA26	CMO
经营活动产生的 现金流量净额增 长率	因子 MTM 的 10 日均 值	26 日指数移动均线	钱德动量摆动指标
KDJ_J	InventoryTRate	CashRateOfSales	RSTR24
随机指标	存货周转率	经营活动产生的现金流 量净额与营业收入之比	24 月相对强势
QuickRatio	MA120	TVMA6	ETP5

速动比率	120 日移动均线	6 日成交金额的移动平均值	5 年平均收益市值比
EMA20	VOL120	ADX	TVSTD6
20 日指数移动均线	120 日平均换手率	平均动向指数	6 日成交金额的标准差
LongDebtToAsset	APBMA	DASREV	FixedAssetsTRate
长期借款与资产总计之比	绝对偏差移动平均	分析师盈收预测变化	固定资产周转率
LongTermDebtToAsset	VEMA10	MA5	LongDebtToWorkingCapital
长期负债与资产总计之比	成交量的 10 日指数移动平均	5 日移动均线	长期负债与营运资金比率
BollUp	PLRC6	BIAS10	VSTD10
上轨线	6 日收盘价格线性回归系数	10 日乖离率	10 日成交量标准差
NonCurrentAssetsRatio	NetProfitGrowRate	TOBT	ASSI
非流动资产比率	净利润增长率	超额流动	对数总资产
Ulcer10	CCI88	EquityTRate	OperatingProfitGrowRate
10 向下的波动性	88 日顺势指标	股东权益周转率	营业利润增长率
TVMA20	SFY12P	ADTM	TotalProfitGrowRate
20 日成交金额的移动平均值	分析师营收预测	动态买卖气指标	利润总额增长率
EquityToAsset	EquityFixedAssetRatio	HBETA	AD20
股东权益比率	股东权益与固定资产比率	历史贝塔	累积/派发线
VOSC	RVI	ATR14	
成交量震荡	相对离散指数	14 日均幅指标	

致谢

本论文是在导师彭虎锋老师的悉心指导下完成的。从入学开始，彭老师就给我们讲解了本次论文研究的相关知识，在老师的带领下，我们逐渐的入门。彭老师的学识渊博以及严谨的治学态度令我们印象深刻。

在生活中，彭老师对我们也是非常关心，尤其是毕业求职阶段，老师给了我巨大的帮助。在此，谨向导师表示崇高的敬意和衷心的感谢！

伴随着毕业论文的完成，两年研究生生活也临近结尾，回望两年来的求学生涯，滋味万千，有即将走向社会的兴奋，也有从此不在是在校学生的感慨。