

The Impact of Macroeconomic Factors on Stock Returns:

Using the S&P 500 Index as An Example

by

Yangui Guo

A Dissertation Presented in Partial Fulfillment  
of the Requirements for the Degree  
Doctor of Business Administration

Approved on March 2024 by the  
Graduate Supervisory Committee:

Xu Wu, Co-Chair  
Hong Yan, Co-Chair  
Huibing Zhang

ARIZONA STATE UNIVERSITY

May 2024

宏观经济因子对股票收益率的影响研究—以标准普尔 500 指数为例

郭延贵

全球金融工商管理博士  
学位论文

研究生管理委员会  
批准于二零二四年三月

吴旭，联席主席  
严弘，联席主席  
张慧冰

亚利桑那州立大学

二零二四年五月

## ABSTRACT

This article uses the S&P 500 index as an example to analyze the impact of macroeconomic factors on stock returns. By using the S&P 500 index data from 1968 to 2020 as the dependent variable, and the monthly data of 221 macroeconomic variables such as the consumer price index and the US mid-term election as the independent variable, this paper finds that: (1) a wavelet denoising method helps to capture the low-frequency and long-term fluctuations in monthly returns of the index, which can effectively remove the short-term fluctuations in returns, better reflect the macroeconomic trend, and improve the power of out-of-sample forecasting. (2) the Granger causality test may be used to pick the top 30 most significant variables, which can be incorporated into several prediction models. Among all the prediction models, the combined prediction algorithm has the best out-of-sample prediction effect. (3) investors need to consider investment practices under timing strategies. Elastic network, scaling principal component analysis, combination prediction, and other algorithms are used to select the time, and the best results are obtained based on the scaling principal component analysis algorithm and the combination prediction algorithm when the transaction fee is set to 5‰. The returns based on these two algorithms have reached 14.00% and 12.59%, their Sharpe ratios are the highest among all algorithms, reaching 0.69 and 0.62, respectively, and this result is significantly better than the historical mean model used as a measurement benchmark (with the average return of 8.14%, and a

Sharpe ratio of 0.34). (4) explore investment practice under a stock-picking strategy. We use methods such as sector rotation strategy and mean-variance of sectors for stock selection, and find that the strategy returns achieved by investing in stocks using the sector rotation strategy are the best, reaching 14.30%, and the Sharpe ratio is the highest at 0.79, significantly better than the benchmark S&P 500 index (with the average return of 8.8% and a Sharpe ratio of 0.57).

## 摘要

本文以标普 500 指数为例，分析宏观经济因子对股票收益率的影响。具体而言，本文使用 1968-2020 年标普 500 指数数据作为因变量，使用居民消费价格指数、美国中期选举等 221 个宏观经济变量的月度数据作为自变量，发现：（1）通过使用小波去噪方法捕捉月度收益率曲线中的低频、长周期波动，能够有效去除收益率曲线中的噪音，使其反映平缓变化的宏观经济趋势，提高样本外预测效果。（2）使用格兰杰因果检验筛选出每期所有自变量中最为显著的 30 个变量投入不同预测模型中进行训练，6 种预测模型中，组合预测的算法样本外预测效果最好。（3）考虑择时策略下的投资实践。使用弹性网络、缩放主成分分析、组合预测等算法进行择时选择，在交易手续费设置为 5‰ 的情况下，基于缩放主成分分析算法和组合预测算法得到的结果最佳。各项指标中，这两种算法的策略收益率分别达到 14.00% 和 12.59%，夏普比率在所有算法中最高，分别达到 0.69 和 0.62，效果均显著优于作为衡量基准的历史均值模型（策略收益率 8.14%，夏普比率 0.34）。（4）考虑选股策略下的投资实践。使用行业轮动策略、行业均值方差等方法进行股票选择，发现使用行业轮动策略进行选股投资得到的策略效果最佳。各项指标中，策略收益率达到 14.30%，夏普比率在所有算法中最高，达到 0.79，效果显著优于作为衡量基准的标普 500 指数（策略收益率 8.8%，夏普比率 0.57）。

## 目录

	页码
表格列表 .....	vii
图表列表 .....	viii
章节	
一、 研究意义 .....	1
二、 文献综述 .....	4
三、 计量经济学方法 .....	8
3.1 小波分析 .....	8
3.2 格兰杰因果关系检验与变量筛选 .....	12
3.3 预测模型 .....	14
3.3.1 最小绝对收缩和选择算法 (LASSO) .....	15
3.3.2 弹性网络算法 (Enet) .....	16
3.3.3 主成分分析算法 (PCA) .....	16
3.3.4 缩放主成分分析算法 (SPCA) .....	18
3.3.5 偏最小二乘法 (PLS) .....	18
3.3.6 历史均值法 (hmean) .....	20
3.3.7 组合预测算法 .....	21
四、 预测结果 .....	22
4.1 数据 .....	22
4.1.1 自变量介绍 .....	22

章节	页码
4.1.2 自变量描述性统计.....	39
4.2 检验指标.....	51
4.2.1 样本外 R 方统计量.....	52
4.2.2 CW-t 检验.....	52
4.3 样本外预测结果.....	53
4.3.1 未小波去噪.....	53
4.3.2 小波去噪后.....	55
4.3.3 不同预测期限对预测效果的影响.....	62
五、投资组合实践.....	65
5.1 衡量指标.....	65
5.1.1 年化收益率 (Annualized rate of return) .....	65
5.1.2 确定性等价收益 (Certainty equivalent return) .....	65
5.1.3 夏普比率 (Sharpe ratio) .....	66
5.1.4 最大回撤 (Max drawdown) .....	66
5.2 手续费计算.....	67
5.3 择时策略.....	67
5.3.1 策略介绍.....	67
5.3.2 策略回测结果.....	68
5.3.3 变量类型对回测结果的影响.....	74

章节	页码
5.3.4 累积财富积累曲线 .....	76
六、选股策略 .....	79
6.1 行业组合 .....	79
6.2 策略介绍 .....	80
6.3 行业样本外预测结果 .....	81
6.4 策略回测结果 .....	86
6.5 累积财富积累曲线 .....	88
七、预测框架与方法的优缺点分析 .....	92
参考文献 .....	95

## 表格列表

表格	页码
1 自变量统计结果 .....	44
2 小波去噪前后样本外预测效果.....	54
3 出现频率最高的 30 个变量介绍 .....	57
4 基于各预测模型的投资组合策略收益 .....	69
5 考虑交易手续费和冲击成本后的投资组合收益情况 .....	71
6 使用不同类型指标的策略回测结果 .....	74
7 行业划分介绍 .....	79
8 各行业样本外预测结果.....	82
9 选股策略投资组合收益情况 .....	86

## 图表列表

图表	页码
1 小波去噪处理后的标普 500 指数超额收益.....	11
2 变量自回归结果 .....	59
3 不同预测期限对预测效果的影响.....	63
4 择时策略的收益率曲线.....	77
5 选股策略的收益率曲线.....	89

## 一、研究意义

资本市场收益率的预测一直是现代金融学的核心问题之一。投资者为了最大化投资收益需要实时预测资本市场收益率；学术研究者则为了提出更好的资产定价模型，深入理解资产收益预测的本质，也需要进行相关研究。在经济条件不稳定、市场波动频繁的背景下，如何选择合适的预测方法和预测因子进行资本市场收益率预测并分析其可预测性的来源成为近年来资产定价领域的研究热点。

本文重点关注宏观经济变量是否能够对美国标准普尔 500 指数（S&P 500 Index）收益率进行有效预测。投资者进行投资时，通常分为选股和择时两种流派。选股投资者进行横截面预测，以预测不同公司股票收益率为目标；择时投资者进行时间序列预测，以标准普尔 500 指数收益率为目标。标准普尔 500 指数是记录美国 500 家上市公司的一个股票指数，由标准普尔公司创建并维护。该指数采样面广、代表性强、精确度高、连续性好，被普遍认为是股票指数期货合约的理想标的，并被学术界广泛使用作为预测目标。

与单只股票相比，标准普尔 500 指数同时包含许多公司的股票，风险更为分散，能够反映更广泛的市场变化。市场参与者密切关注该指数，因为它的表现代表了美国的宏观经济状况。标准普尔 500 指数的变动主要对系统性风险做出反应，能够有效减少非系统性风险。例如，从 1969 年到 1981 年，标准普尔 500 指数逐渐下降，同时美国经历了增长停滞和高通货膨胀。在 2008 年金融危机和大衰退期间，标准普尔 500 指数下跌了 46.13%。在 2020 年，冠状病毒大流行使世界经济陷入普遍衰退，股票市场陷入困境，标准普尔 500 指数暴跌近 35%。标普 500 指数对于投资者来说意义重大。它是一个重要的经济表现指标，可以反映出美国经济的整体状况和发展趋势。由于指数包含来自各个行业

的公司，因此不会受某个行业或公司的单一事件的影响，而是反映整体市场的走势。因此，投资者可以使用标普 500 指数来评估他们的投资组合在整个市场中的表现。

此外，标普 500 指数也是投资者的重要投资和分析工具之一。准确预测标普 500 指数的走势对投资者非常有利。如果投资者能够准确预测标普 500 指数的走势，就能够制定更合理和准确的投资策略，从而提高投资收益。

因此，标普 500 指数的重要性不可忽视。考虑到标普 500 指数相关资料的完整、透明、全球参与等特征，相较于港股，沪深，日经指数等指数，标普 500 指数是更为优质的预测目标。

从理论上讲，股票市场的表现一般应该取决于宏观经济状态（Bansal 和 Yaron, 2004）。在这个大数据时代，存在着广泛的宏观经济指标，通过对其展开研究，可以深入了解经济的整体健康状况。宏观经济因素在预测股票回报方面起着至关重要的作用。这些因素是反映经济整体健康状况的经济指标，并提供了对股票市场可能表现的洞察。宏观经济因素多种多样，包括通货膨胀、利率、国内生产总值（GDP）和失业率等多种因素。这些因素能够影响整体投资环境，进而影响股市。通过分析这些因素，分析师可以深入了解经济方向，进而了解标准普尔 500 指数。

然而，正如 Goyal 和 Welch（2008）所记录的那样，相对于简单的历史平均值，典型经济变量的样本外预测表现并不令人满意。受 Goyal 和 Welch（2008）的启发，许多研究集中在探索股票回报的有效预测因子上。一些很好的例子包括油价变化（Driesprong 等人，2008），方差风险溢价（Bollerslev 等人，2009），政策不确定性（Brogaard 和 Detzel，2015），一致的投资者情绪（Huang 等人，2015），空头利息（Rapach 等人，

2016), 新闻隐含波动 (Manela 和 Moreira, 2017), 经理情绪 (Jiang 等人, 2019) 和行业等相关 (Wang 等人, 2020)。

与以往的研究一致, 本文着重于构建驱动股票收益系统变化的不可观察潜在因素的有效代理变量, 通过考虑股票回报动态频域的新方法实现这一目标。在利用机器学习算法进行财务预测的相关研究中 (例如, Huang et al., 2015; Lin, 2018; Zhang et al., 2019; Huang et al., 2021), 目标变量通常设置为资产回报率, 而本文的变量则是逐渐变化的宏观经济指标。因此, 由于其剧烈的波动, 已实现的月度股票收益显然不是一个合适的选择。为了解决这个问题, 本文对月回报序列进行小波频域分解, 并使用其长期 (低频) 成分作为宏观经济状态的代理。主要目的是消除公认的宏观经济变量的噪声成分, 并提取反映实际经济基本面的单个变量的共同预测成分。为了实现这一点, 本文采用了小波去噪方法, 利用小波变换对信号进行去噪处理, 捕捉收益率曲线中超过 8 个月的长期趋势, 排除中高频率的噪声影响, 并使用最小绝对收缩和选择算法、弹性网络算法、主成分分析法、缩放主成分分析、偏最小二乘法和组合方法等方法对标普 500 指数超额收益率进行预测。基于预测结果, 本文分别使用择时策略与行业选股策略, 进行指数的买入与卖出操作, 并根据收益率、确定性等价收益、夏普比率、最大回撤等指标衡量投资策略的效果。

## 二、文献综述

John H. Cochrane (2008) 指出, 对于股票市场收益的预测在金融学 研究领域具有重要地位。但是股票市场具有独特的低信号、高噪音、具有时变性的特点, 对其进行预测非常困难, 目前的理论与实证研究都处于起步阶段。其次, 投资者为了获取超额收益, 使用前沿预测模型与预测因子, 在激烈的竞争下, 资本市场朝着减弱可预测性的方向动态调整, 因此而进一步提高了收益预测的难度。Zhou (2013, 2017) 在理论上证明了月度股票收益的可预测性 $R^2$ 有一个上界, 这个上界是 8%, 然而这个上界仍然被学术界和业界认为是非常松弛的。目前最前沿的预测模型和预测因子所能提供的预测能力仍然非常低, 如何准确高效地对资本市场收益进行预测是一个有趣而困难的研究领域。

在金融预测领域, 宏观变量被广泛使用: Campbell (2000) 提出标准普尔 500 收益率与利率利差和通货膨胀等因素有关。Shiu-Sheng Chen (2009) 使用宏观变量来预测股票市场的衰退, 认为收益率曲线利差和通胀率是用来预测美国股市经济衰退最优指标。该研究比较了宏观经济变量对股市经济衰退预测与股票回报可预测性的预测效果, 认为使用宏观经济变量对股市经济衰退进行预测可以获得更好的效果。Kofi O. Nti (2019) 基于拥有特殊特征选择规则的随机森林算法和长短期记忆循环神经网络算法 (LSTM-RNN), 使用宏观变量对资本市场的收益进行预测研究。该研究强调了宏观经济变量与不同行业股票价格变动存在关联性上的差异, 宏观变量在预测不同行业的股票价格变动时具有不同的权重。Chowdhury, Shah & Abu (2006) 探索了宏观变量的波动率与股票市场波动率的关系。研究发现宏观变量的波动率作为市场风险的来源, 与股票市场波动率以及风险调整后的预期收益呈正相关关系。Ludvigson (2007) 提出使用动态因子分析的方法从大量经济

变量中提取主要的公共宏观因子，使用少量宏观因子有效地总结了经济时间序列中的大量信息，对股票市场超额回报的预测效果良好。

与此同时，许多新的宏观预测指标被发现和提出：Jedrzej, Katrin & Wisniewski (2007) 提出使用全国性选举作为宏观因子，预测股票市场的波动情况。研究发现，在全国性选举日前后一周的窗口中，该国股票市场的波动率会增加到平时的两倍。强制性投票法案通过提高选民投票率与选举前调查的准确性的方式，在一定程度上减少了不确定性，降低了股票市场波动率。Jacobsen, Marshall & Visaltanachoti (2019) 发现工业原材料价格可以被认为是经济和股票市场的重要领先指标，标准普尔 500 指数收益率的变动与工业原材料价格有关，并与当时经济所处发展状况存在一定关联。工业原材料的价格在经济扩张期的上涨将带来标准普尔 500 指数收益率的相应上涨，而在经济衰退期则可能导致标准普尔 500 收益率的下跌。Afees A.Salisu (2020) 将新冠疫情相关的健康新闻纳入股票回报预测模型，取得了良好的预测效果。研究发现健康新闻因子对股票收益具有负面影响，而健康新闻搜索量的增加会加剧其对股票收益的负面影响。

宏观变量在股票收益预测领域占据着十分重要的地位，国内外学者已在这一领域进行了许多研究。然而，以上的研究主要集中于在样本内使用宏观变量来解释标准普尔 500 指数的变动，Goyal & Welch (2008) 研究发现众多具有样本内解释能力的预测变量，并不能提高样本外预测的精度。再考虑到宏观数据从产生到被获取的过程中，存在很长的延迟时间，因而使用宏观变量对标准普尔 500 指数进行准确预测显然要更加困难。

在预测模型的选用方面，机器学习算法被广泛应用于大数据的分析当中 Han 等人 (2019) 提出使用超过一百个特征对横截面收益率进行预测的流程。Wold (1966, 1975)

首先提出了 PLS 预测方法，并由 Kelly 和 Pruitt（2013, 2015）进行进一步扩展。与广泛使用的主成分（PC）分析不同，它在从自变量中提取公因子时利用了因变量，因此 PLS 因子与因变量的相关性更高，比 PC 因子更具有预测能力。Rapach 等人（2013）最先在金融领域应用最小绝对收缩算法，他们分析了高维度环境下的国际股票收益率之间的领先-滞后关系。Gu 等人（2019）使用了包括最小绝对收缩算法在内的机器学习算法，用于分析时间序列月度个股回报率的不可预测性。Chinco 等人（2019）使用最小绝对收缩算法以一分钟的领先时长预测个股回报。Freyberger 等人（2017）使用 LASSO 的非参数方法来测试哪些特征能够提供预期收益横截面的信息，拟合特征值与股票预期收益率之间的非线性关系。与线性面板回归方法相比，方法具有更高的样本外解释力，将夏普比率提高了 50%。Kozak 等人（2020）在贝叶斯背景中使用最小绝对收缩算法对随机贴现因子进行建模。

Rapach, Strauss & Zhou（2009）认为单个回归模型具有不确定性和不稳定性，严重影响了模型的预测能力，建议使用多个预测模型的组合以整合来自多个经济变量的信息并显著降低预测波动。Green, Hand, and Zhang（2017）使用 94 个公司特征，通过普通最小二乘法与加权最小二乘法对美股股票收益率进行预测。研究发现大多数公司在每个时间点平均大约 30 个公司特征具有显著重要性。Han 等人（2019）使用最小绝对收缩和选择算法与弹性网络等机器学习算法，在避免过度拟合的情况下提取大量经济变量中的有效信息进行池化预测。Hong, Jiang & Meng（2022）整合了最小绝对收缩和选择算法弹性网络算法、最小二乘法、偏最小二乘法、缩放主成分分析算法、随机森林算法，采用以上算法的组合形式展开预测。研究针对来自 80 万篇《华尔街日报》文章的全文内容进

行定量，以通货膨胀率作为预测指标。结果显示组合方法在样本内与样本外的预测效果均优于基准模型。

在构建选股投资策略方面：Markowitz（1952）提出了在静态环境下在风险资产之间分配财富的最佳规则，也即投资者只关心投资组合回报的均值和方差。但由于该方法容易导致投资组合中的某些权重过高，从而随着时间的推移，产生巨大波动且样本外表现不佳，处理估计误差的问题一直是该模型有待解决的问题之一。贝叶斯方法在构建投资组合的过程中时常被用来进行误差估计，其多种实现范围从纯粹依赖于弥散先验的统计方法（Barry, 1974 年；Bawa, Brown & Klein, 1979），到“收缩估计量”（Jobson, Korkie 和 Ratti, 1979；Jobson 和 Korkie, 1980 年；Jorion, 1990），到最近的方法依赖于资产定价模型来建立先验（Pastor, 2000 年；Pastor 和 Stambaugh, 2000）。同样常用的是使用非贝叶斯方法来估计误差，其中包括“稳健”的投资组合分配规则（Goldfarb 和 Iyengar, 2003 年；Garlappi, Uppal 和 Wang, 2007 年）；投资组合规则设计最佳分散市场和估计风险（Kan 和 Zhou, 2007 年）；利用因子结构强加的矩限制的投资组合回报率（MacKinlay 和 Pastor, 2000 年）；专注于减少的方法估计协方差矩阵的错误（Best 和 Grauer, 1992 年；Chan, Karceski 和 Lakonishok, 1999 年；Ledoit 和 Wolf, 2004）；最后，施加卖空约束的投资组合规则（Frost 和 Savarino, 1988 年；Chopra, 2013；Jagannathan 和 Ma, 2003）。DeMiguel, Garlappi 和 Uppal（2009）以 1/N 策略为基准，研究比较了 14 种最佳资产配置模型的表现。这 14 中资产配置模型包括基于样本均值方差策略、贝叶斯漫反射先验模型、最小方差模型等。研究发现各种优化策略下，没有一个模型可以始终如一地提供夏普比率或确定性等价收益高于 1/N 投资组合，即收益率较低。

### 三、 计量经济学方法

#### 3.1 小波分析

小波分解是一种有效的信号分解技术，可以应用于金融时间序列分析。金融时间序列是指按时间顺序排列的数据，例如股票价格和汇率等。通过对金融时间序列进行小波分解，可以将其分解成不同时间尺度和频率的组成部分，这些部分可用于分析时间序列的趋势、周期性等特征，从而更好地预测未来的趋势和风险。

小波变换多分辨率分析（DWT-MRA）是小波分解的延伸，利用小波基函数的正交性和尺度变换特性，将信号分解成多个尺度的子信号，这些子信号分别代表信号的不同频率成分，实现了信号的多分辨率分析。在分解过程中，每一层分解的子信号都可以进一步分解，形成一个小波分解树，树的深度代表信号的分辨率等级，也称为分解层数。每一层分解的子信号可代表不同频率的信号成分，越靠近根节点的子信号代表越低频的信号成分。

小波变换多分辨率分析（DWT-MRA）通常用作信号处理方法，在有限的时间窗口内分解目标序列。通过对信号进行小波分解，DWT-MRA 可将其分解成不同频率的子信号，并对这些子信号进行多尺度分析。该方法可以有效提取信号的局部特征，实现信号的时频局部化处理。分解的大小取决于过滤器窗口的数量。具体来说，当选择较小的滤波器窗口来捕获更高频率的特征时，将分析和生成原始序列中更详细的分量，即不同的窗口大小对应于不同的多分辨率（频率）分量的子序列。对于时间序列  $y_t$ ，多分辨率表示由以下公式给出：

$$y_t = y_t^{S_j} + y_t^{D_j} + y_t^{D_{j-1}} + \dots + y_t^{D_1}$$

其中  $y_t^{S_j}$  是平滑分量， $y_t^{D_j}$  是详图分量， $j = 1, \dots, j$ ， $j$  的值越大， $y_t^{D_j}$  的信息频率越高。

在金融领域，预测股票回报是一项重要的任务。本文旨在通过研究收益率曲线来分析宏观经济状态和预测市场趋势。然而，月度股票收益率曲线常常存在剧烈波动和噪音，这为研究带来了很大挑战。为了解决这个问题，本文采用小波频域分解技术对收益率曲线进行分解，以提取不同频率和时间尺度上的成分。具体地，本文使用小波分解将月度收益序列分解成长期（低频）、中期和短期成分。长期成分代表了宏观经济状态的平缓变化，可作为反映经济基本面的目标变量；而中短期成分则可用于预测市场的波动和趋势。在此背景下，本文借鉴 Faria 和 Verona（2018、2020、2021）的研究，使用最大重叠离散小波变换（MODWT）和反射边界条件的 Haar 小波滤波器对回报序列进行分解，得到三个分量。其中，平滑分量表示低频（长期）行为，细节分量 1 和细节分量 2 代表中高频（中短期）行为。通过这些过程，本文将原始序列表示为来自不同频域的三个子序列的组合。通过小波分解技术，可以有效地去除收益率曲线中的噪音和波动性，提取出长期趋势和宏观经济状态的信息。同时，多尺度特性可以将收益率曲线分解成不同时间尺度的成分，挖掘市场的多种特征和规律。这些成分可以作为重要的预测因子，用于预测市场的未来走势和趋势。

为了避免前瞻偏差，本文仅在预测期间  $t + 1$  的回报时，将回报序列分解到周期  $t$ 。这样可以确保使用的信息只来自  $t$  时期及之前的数据，避免了未来信息的泄漏。使用 MODWT 和 Haar 小波滤波器对股票回报序列进行分解，可以帮助更好地理解序列的行为模式，并提供更准确的预测信息。同时，避免前瞻偏差可以确保使用的信息是可靠的，并提高预测的准确性。

在本文的分析中，应用  $J = 2$  级 MRA:

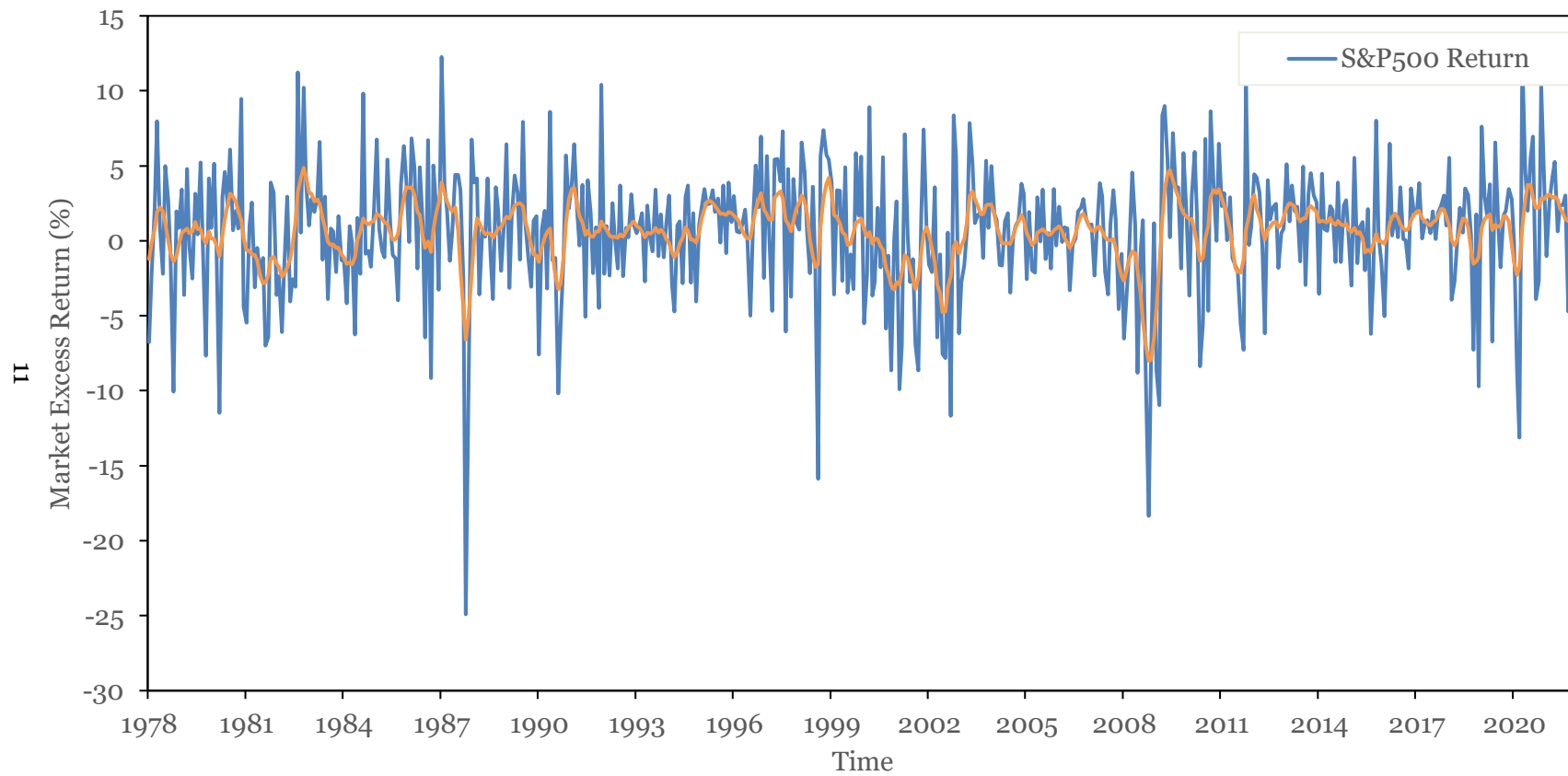
$$y_t = y_t^{S_2} + y_t^{D_1} + y_t^{D_2}$$

其中， $y_t^{D_1}$ 与 $y_t^{D_2}$ 为 2 个小波细节分量， $y_t^{S_2}$ 为小波平滑分量。由于本文在分析中使用了月度数据，第一个小波细节分量 $y_t^{D_1}$ 捕捉 2 到 4 个月之间的振荡，第二个细节分量 $y_t^{D_2}$ 捕获 4 到 8 个月之间的振荡，平滑组件 $y_t^{S_2}$ 捕获周期超过 8 个月的振荡。这意味着，如果存在周期大于 8 个月的经济波动，这些波动将会反映在平滑分量 $y_t^{S_2}$ 中。在本文中，考虑到宏观变量本身自相关性很高，变化非常缓慢，因此假定收益率序列的高频和中频分量与宏观变量无关，选取平滑组件 $y_t^{S_2}$ 的低频时间序列分量进行后续研究。

小波去噪处理前后的标普 500 指数超额收益曲线如图 1 所示：

图 1 小波去噪处理后的标普 500 指数超额收益

S&P500 Return And Denoised S&P500 Return



本文介绍了对标普 500 指数超额收益率曲线进行小波去噪处理的方法及效果。原始的曲线波动剧烈，不易观察到整体趋势特征，因此需要进行去噪处理。本文的研究关注的是反映平缓变化的宏观经济状态，因此处理后的曲线（橙色）更加平滑，波动范围变小，同时仍然保留了原始曲线（蓝色）的整体趋势特征，让人能够更好地理解收益率的整体趋势。图中展示了 1978 年-2020 年间小波去噪处理前后的标普 500 指数超额收益率曲线，横轴表示时间，纵轴表示超额收益率。本文使用小波去噪处理对原始收益率曲线进行了平滑，这有助于更好地观察收益率曲线的趋势特征。本研究采用小波去噪处理方法，能够更好地反映出宏观经济状态的整体趋势，提高了对收益率变化的理解。

### 3.2 格兰杰因果关系检验与变量筛选

在海量数据中，能够有效预测未来股票市场收益的信号总是稀疏的，绝大多数变量是没有预测能力的噪音，而且他们的预测能力也常常是时变的，一个变量常常只在某些时期有效，在其他的时期就失去了预测能力。考虑到单变量回归是最为普遍的变量筛选方法，本文采用格兰杰因果检验的方法，来对变量进行预筛选。

格兰杰因果关系检验是一种统计假设检验，用于确定一个时间序列是否有助于预测另一个时间序列，于 1969 年首次提出。通常，回归被认为是反映“纯粹”的相关性，但克莱夫·格兰杰认为，经济学中的因果关系可以通过测量使用另一个时间序列的先验值来预测一个时间序列的未来值的能力来测试。

如果时间序列  $X$  是可以获得的，那么可以对  $X$  的滞后值（也包括  $Y$  的滞后值）进行一系列  $t$  检验和  $F$  检验，来验证原假设： $X$  不是  $Y$  的格兰杰原因。称时间序列  $X$  是  $Y$  的格兰杰原因，因为这些  $X$  值提供有关  $Y$  未来值的统计显著信息。

格兰杰因果检验的数学表达为：

$$y_t = a_0 + a_1y_{t-1} + a_2y_{t-2} + \dots + a_my_{t-m} + b_px_{t-p} + \dots + b_qx_{t-q}$$

其中， $m$  为自回归滞后长度， $p$  是自变量回归最短的滞后长度， $q$  是最长的滞后长度。股票收益时间序列存在以下两个特点：①相邻两期股票收益相关系数接近 0。②滞后超过两期的  $x_t$  对  $y_t$  没有预测能力。因此本文中设置  $m=0$ ， $p=q=1$ 。Huang, et al. (2022) 中使用了相同的方法来评估不同因子的预测能力，从而对不同因子的方差进行缩放。Kelly, B., & Pruitt, S. (2015) 也使用了相同的方法和参数设置，来评估不同预测变量的预测能力。

表达式为：

$$y_t = a_0 + b_1x_{t-1}$$

基于以上表达式进行回归拟合，对原假设进行检验：

$$H_0: a_0 = b_1 = 0$$

当且仅当回归结果中不保留  $x$  的滞后值时，才接受  $x$  不是  $y$  的格兰杰原因的原假设。

根据回归结果，在每个时期都能得到不同变量斜率  $b_1$  的  $t$  值，然后依据  $t$  值的绝对值对变量进行排序，选择  $t$  值最大的 30 个变量作为特征值，使用机器学习模型进行预测。也就是说，在对每一期收益进行预测以前，本文都使用了格兰杰因果关系检验的方法，基于该时期所对应的样本内数据进行变量筛选。这即意味着，变量筛选的结果在时间维度上是动态的，不同时间节点上筛选出的 30 个最为显著的自变量存在差异。

### 3.3 预测模型

线性回归模型是一种经典的机器学习算法，用于分析自变量和因变量之间的线性关系。当需要预测未来股票市场指数超过无风险利率的收益率时，可以使用该模型来估计收益率的取值。该模型的公式可以表示为：

$$r_{t+1} = \beta_t x_t + \varepsilon_{t+1}$$

其中 $r_{t+1}$ 是股票市场指数超过无风险利率的收益率， $x_t$ 是具有预测能力的变量， $\varepsilon_{t+1}$ 是干扰项。

参考 Welch 和 Goyal (2008)，本文使用递归（扩展）估计窗口进行股票溢价的样本外预测。对股票市场指数超过无风险利率的收益率根据变量进行拆解，变量 $x_m$ 对应预测收益率 $r_{m+1}$ 。具体而言，本文首先将 $r_t$ 和 $x_t$ 的 $T$ 个观测值的总样本分为由前 $m$ 个观测值组成的样本内部分和由最后 $q$ 个观测值组成的样本外部分。基于预测因子 $x_t$ 的股票溢价的初始样本外预测由下式给出：

$$\hat{r}_{m+1} = \hat{\beta}_m x_m$$

其中 $\hat{\beta}_m$ 是在各算法对应计算得到的 $\beta_m$ 的估计值，由基于 $\{x_t\}_{t=1}^{m-1}$ 的回归 $\{r_t\}_{t=2}^m$ 得到，下一个样本外预测由下式给出：

$$\hat{r}_{m+2} = \hat{\beta}_{m+1} x_{m+1}$$

其中 $\hat{\beta}_{m+1}$ 由基于 $\{x_t\}_{t=1}^m$ 的回归 $\{r_t\}_{t=2}^{m+1}$ 得到。

基于 $x_t$ ， $\{r_{t+1}\}_{t=m}^{T-1}$ ，按照以上方式在样本外区间进行依次推导，得到由 $q$ 个观测值组成的样本外股票溢价预测序列。

在本研究中，第一个样本外预测对应的样本内观测值为 1968 年 1 月-1977 年 12 月间的月度宏观变量数据与对应的标普 500 指数收益率的真实值，由此得到 1978 年 1 月的样本外标普 500 指数收益率的预测值；第二个样本外预测对应的样本内观测值为 1968-1978 年 1 月的月度宏观变量数据与对应的标普 500 指数收益率的真实值，由此得到 1978 年 2 月的样本外标普 500 收益率的预测值...以此类推最终得到 1978-2020 年间，标普 500 指数收益率的预测值序列。

### 3.3.1 最小绝对收缩和选择算法（LASSO）

Tibshirani（1996）提出的最小绝对收缩和选择算法本质上是一个惩罚线性模型，是一种常用的统计学习方法，它可以在高维数据集上进行特征选择和模型建立。与传统的最小二乘回归方法相比，LASSO 可以在保持预测准确度的前提下，同时降低模型复杂度和提高可解释性。

LASSO 的核心思想是在目标函数中添加一个 L1 正则化项，这个正则化项的系数可以控制模型的复杂度和特征选择的效果。在 L1 正则化项的影响下，某些特征的系数会被缩减至零，从而实现特征选择。LASSO 的优点在于它能够自动选择特征，不需要额外的特征工程，同时还能够抵抗数据中的噪声和过拟合。

基于预测因子  $x_m$  的股票溢价的初始样本外预测由下式给出：

$$\hat{r}_{m+1} = \hat{\beta}_m x_m$$

其中  $\hat{\beta}_m$  可以由下式得到：

$$\hat{\beta}_m = \arg \min_{\beta} \left( \sum_{t=1}^m (r_t - \beta' x_{t-1})^2 + \lambda \sum_{k=1}^K |\beta_k| \right)$$

下一个样本外预测由下式给出：

$$\hat{r}_{m+2} = \hat{\beta}_{m+1} x_{m+1}$$

其中  $\hat{\beta}_{m+1}$  由基于  $\{x_t\}_{t=1}^m$  的回归  $\{r_t\}_{t=2}^{m+1}$  得到。

### 3.3.2 弹性网络算法 (Enet)

弹性网络是 Zou 和 Hastie (2005) 引入的另一个惩罚线性模型。弹性网络是一种线性回归的扩展，旨在克服传统线性回归中存在的过拟合问题。与传统的线性回归不同，弹性网络可以同时进行特征选择和参数估计。弹性网络使用 L1 和 L2 正则化项来控制模型的复杂度和泛化能力。L1 正则化项可以将一些系数缩小到 0，从而实现特征选择的目的，因此可以减少一些不重要的特征对模型的影响。而 L2 正则化项则可以缩小所有系数，从而减少模型的过拟合风险，提高模型的泛化能力，克服了最小绝对收缩和选择算法的一些限制。

其样本外预测由下式给出：

$$\hat{r}_{m+1} = \hat{\beta}_m x_m,$$

其中  $\hat{\beta}_m$  可以通过以下方式估算：

$$\hat{\beta}_m = \arg \min_{\beta} \left( \frac{\sum_{t=1}^m (r_t - \beta' x_{t-1})^2}{m} + \alpha \lambda \sum_{k=1}^K \beta_k^2 + \frac{(1-\alpha)}{2} \lambda \sum_{k=1}^K |\beta_k| \right)$$

下一个样本外预测由下式给出：

$$\hat{r}_{m+2} = \hat{\beta}_{m+1} x_{m+1}$$

其中  $\hat{\beta}_{m+1}$  由基于  $\{x_t\}_{t=1}^m$  的回归  $\{r_t\}_{t=2}^{m+1}$  得到。

### 3.3.3 主成分分析算法 (PCA)

主成分分析法是宏观经济和金融预测中非常流行的因子模型，是一种常用的线性降维技术。其主要目的是通过找到数据中的主要成分，从而减少数据的维度。

主成分分析的基本思想是通过线性变换将原始数据映射到一个新的空间中，使得新的空间中的数据尽可能地保留原始数据的信息，并且不同特征之间的相关性尽可能地降低。具体来说，主成分分析将原始数据转换为一组新的特征，这些新的特征被称为主成分，每个主成分是原始特征的线性组合。通过保留前几个主成分，可以达到降低数据维度的目的，同时保留了数据的主要特征。考虑到第一个主成分因子在很大程度上已经解释了市场变化，并且使用多个主成分因子会显著增加计算时间，降低预测效率，本文中仅保留第一个主成分。在主成分分析法中，若干变量被转化为正交分量，前几个主要分量倾向于解释所有变量中方差最大的，分量等级越高，它能够解释的方差就越大。

假设所有预测变量都由一些常见的潜在因素驱动，

$$\mathbf{x}_t = \Lambda \mathbf{f}_t + \mathbf{e}_t$$

其中 $\mathbf{f}_t$ 是来自 $\mathbf{x}_t$ 的主分量的  $n$  向量， $n$  远小于预测变量的数量。 $\mathbf{f}_t$ 可以通过最大化预测变量的方差来估计。由此得到的预测模型为：

$$\hat{r}_{m+1} = \hat{\beta}_m \mathbf{f}_m$$

其中 $\hat{\beta}_m$ 是 $\beta_m$ 的普通最小二乘（OLS）估计值，由基于 $\{\mathbf{f}_t\}_{t=1}^m$ 的回归 $\{r_t\}_{t=2}^m$ 得到。

然而，PCA 也存在一些局限性，这些局限性可能导致其在某些情况下无法捕获原始预测因子中包含的有用信息。由于主成分分析法属于无监督学习方法，弊端在于其完全忽略目标信息，即股票收益率。由于 PCA 只关注数据的方差和协方差，它无法捕获股票收益率与其他预测因子之间的非线性关系，从而限制了它的预测能力。因此 PCA 方法不一定能够捕获原始预测因子中包含的预测股票收益率的有用信息。另一个局限性是 PCA 在处理非线性数据时可能会出现的问题。由于 PCA 基于数据的方差和协方差，它可能无法捕获

非线性数据中的主要模式，这可能会导致 PCA 在预测中失效。

### 3.3.4 缩放主成分分析算法（SPCA）

缩放主成分分析算法旨在解决主成分分析法中存在的容易忽略目标信息的问题。与偏最小二乘法不同，缩放主成分分析法直接为那些具有更强预测能力的预测因子分配更多的权重。SPCA 在求解主成分时引入了一个额外的缩放因子，从而使得每个主成分都具有不同的方差，这有助于减少数据的噪声影响，并提高数据的解释能力。模型的表达式为：

$$y_{m+h} = v_i + \gamma_i X_{i,m} + u_{i,m+h}, i = 1, \dots, N$$

通过回归得到每一个预测因子对应的系数  $\hat{\gamma}_i$ ，进而得到缩放后的自变量  $(\hat{\gamma}_1 X_{1,m}, \hat{\gamma}_2 X_{2,m}, \dots, \hat{\gamma}_N X_{N,m})$ ，使用主成分分析法对缩放后的自变量提取主成分。具体而言，计算  $T \times T$  矩阵  $M_{XX}^\circ = \frac{1}{N} \sum_{i=1}^N \hat{\gamma}_i \dot{X}_i (\hat{\gamma}_i \dot{X}_i)'$ ，其中  $\dot{X}_i$  代表自变量  $i$  的退化矢量，也即减去均值后的变量值。缩放主成分  $\widehat{F}^{SPCA}$  等于矩阵  $M_{XX}^\circ$  的特征向量的  $\sqrt{T}$  倍，其中特征向量为对应特征根从大到小排序的前  $r$  个。本文中  $h$  取 1，由此得到预测模型：

$$\widehat{y_{m+1}}^{SPCA} = \hat{\beta}_m \widehat{F}^{SPCA}_m$$

其中， $\hat{\beta}_m$  是  $\beta_m$  的普通最小二乘（OLS）估计值，由基于  $\{\widehat{F}^{SPCA}_t\}_{t=1}^{m-1}$  的回归  $\{r_t\}_{t=2}^m$  得到。

本文中  $r$  取 1，也即仅保留第一个主成分。

### 3.3.5 偏最小二乘法（PLS）

偏最小二乘法结合了主成分分析和多元回归分析的优点，能够处理高维数据集中存在的共线性和噪声等问题。主成分分析法的一个弱点是它在构建时忽略了目标信息因素，导致预测能力差。Will（1966）提出的偏最小二乘法方法通过最大化自变量和因变量之间的协方差，将自变量和因变量投影到新空间来找到线性回归模型，来克服这一困难。假定宏

观经济变量能够通过线性模型解释未来一期的预期超额股票收益率：

$$E_t(R_{t+1}) = \alpha + \beta S_t$$

在使用宏观经济变量预测资产回报率的情况下， $S_t$ 可以被视为宏观经济因素的一个综合体，尽管它是真实的，但不可直接观测。这些宏观经济变量对于预测资产回报率至关重要。因此，已实现的股票回报率等于在这些宏观经济变量影响下的条件预期回报率加上不可预测的冲击：

$$R_{t+1} = E_t(R_{t+1}) + \varepsilon_{t+1}$$

其中 $\varepsilon_{t+1}$ 是不可预测的，且与 $S_t$ 无关。

令 $x_t = (x_{1,t}, \dots, x_{N,t})'$ 代表第 $t$ 期( $t = 1, \dots, T$ )宏观经济变量的 $N \times 1$ 向量，假设 $x_{i,t}$  ( $i = 1, \dots, N$ )具有因子结构：

$$x_{i,t} = \eta_{i,0} + \eta_{i,1}S_t + \eta_{i,2}E_t + e_{i,t}, \quad i = 1, \dots, N$$

其中， $S_t$ 是对预测资产收益率至关重要的宏观经济因素， $\eta_{i,1}$ 是概括宏观经济变量 $x_{i,1}$ 对 $S_t$ 变动敏感度的因子载荷， $E_t$ 是与收益率无关的所有变量的共同近似误差分量， $e_{i,t}$ 是仅与指标 $i$ 相关的特异性噪声。

偏最小二乘法可以通过以下两步来实现。第一步，进行 $N$ 个时间序列回归。也就是说，对于每个宏观经济变量 $x_i$ ，将 $x_{i,t-1}$ 与常数和已实现股票回报率 $R_t$ 进行时间序列回归：

$$x_{i,t-1} = \pi_{i,0} + \pi_i R_t + u_{i,t-1}, \quad t = 1, \dots, T$$

载荷 $\pi_i$ 反映了每个宏观经济因素 $x_{i,t-1}$ 宏观经济因素 $S_{t-1}$ 的敏感性。由于 $R_t$ 的预期成分由 $S_{t-1}$ 驱动，因此情绪代理与预期股票回报相关，与不可预测的回报冲击不相关。因此，第一阶段时间序列中的系数 $\pi_i$ 近似地描述了每个宏观经济变量如何依赖于真实的宏观经济因素。

第二步，进行  $T$  跨期回归。更具体地说，对于每个时间段  $t$ ，基于相应载荷  $\hat{\pi}_i$  进行关于  $x_{i,t}$  的横截面回归：

$$x_{i,t} = c_t + S_t^{PLS} \hat{\pi}_i + v_{i,t}, \quad i = 1, \dots, N$$

其中， $S_t^{PLS}$  即上式中的斜率，是估计的宏观经济因素的综合体。由此可以将估计得到的  $S_t^{PLS}$  代入线性模型计算得到对未来收益率的预测值。

### 3.3.6 历史均值法 (hmean)

股票收益率历史均值法是一种常见的投资分析方法，它利用股票过去的收益率数据来预测未来的收益率。这个方法的基本思想是通过对股票收益率的历史数据进行分析，计算出历史收益率的平均值，并将这个平均值作为未来收益率的预测值。

在进行预测时，投资者将股票历史平均收益率作为股票未来收益率的预测值，这样就可以更准确地预测股票的未來表现。作为一个恒定的预期股票收益率，股票溢价的历史平均值计算公式如下：

$$\bar{r}_{t+1} = \frac{1}{t} \sum_{j=1}^t r_j$$

这被称作历史均值法基准预测模型。直观而言，如果  $x_t$  包含对预测股票溢价有用的信息，那么  $\hat{r}_{t+1}$  应该比  $\bar{r}_{t+1}$  表现更好。

然而需要注意的是，该方法存在一定的局限性，因为它不能预测股票市场中的异常情况，如经济危机、政治风险等因素对股票价格的影响。因此，投资者需要综合考虑其他因素，如公司业绩、市场前景、政策变化等，以做出更加准确的投资决策。

在本文中，将历史均值模型设定为基准模型，将机器学习模型的预测效果与历史均值模型作比较，以衡量与评估预测模型的准确性。

### 3.3.7 组合预测算法

组合预测方法是一种常见的预测技术，它基于多个不同的预测模型来生成更准确的预测结果。该方法可以利用多个模型之间的差异性，将它们结合在一起，以获得更为准确的预测结果。组合预测方法可以利用多个模型的优点，更好地捕捉数据的潜在关系，提高预测准确率。当组合预测方法中的多个模型产生不同的预测结果时，可以通过统计方法来综合这些结果，得到更为准确的预测结果。组合预测方法可以降低单一模型因为噪声和数据不足等原因而导致的预测误差。通过组合多个模型的预测结果，可以减少误差和偏差，提高整体的鲁棒性。此外，组合预测方法还可以增强模型的解释性。通过使用多个模型进行预测，可以更好地理解模型是如何做出预测的，并且能够更好地理解不同模型的偏差和优点。

具体而言，在时间  $t$  处对  $r_{t+1}$  进行的组合预测  $\hat{r}_{t+1}$  由分别使用  $N$  种预测方法所得到的预测值  $\hat{r}_{i,t+1}$  取加权平均得到。

$$\hat{r}_{t+1} = \sum_{i=1}^N \omega_{i,t} \hat{r}_{i,t+1}$$

其中  $\omega_{i,t}$  是在时间  $t$  所估计的第  $i$  种预测方法的事前组合权重。在本文中，组合预测方法为等权重组合， $N=5$ ，5 种预测方法分别为最小绝对收缩和选择算法、弹性网络算法、主成分分析算法、缩放主成分分析算法和偏最小二乘算法，权重均为 20%。

## 四、预测结果

### 4.1 数据

#### 4.1.1 自变量介绍

本文原始数据中共包含 221 个宏观数据指标，包括 15 个经济指标、7 个市场情绪指标、45 个技术指标、143 个宏观经济指标和 11 个政治选举指标，以下对变量进行介绍：

##### i. 经济指标（15 个）

- 1) **DP**: 股息价格比，等于股息的对数与股票价格的对数的差值。
- 2) **DY**: 股息率，等于股息的对数与滞后股票价格的对数的差值。
- 3) **EP**: 等于股票盈余的对数与股票价格的对数的差值。
- 4) **VOL**: 股票波动率，使用过去 12 个月的月频股票收益率计算得到。
- 5) **TBL**: 国库券利率，1920 年至 1933 年的国库券利率来自于 NBER 的财政部宏观历史数据库的三个月美国储蓄债券利率。从 1934 年到 2003 年的国库券利率来自 3 个月的国库券的二级市场利率，数据从经济圣路易斯联邦储备银行（FRED）的研究数据库中得到。
- 6) **BILL**: TBL 减去 TBL 过去 12 个月的均值。
- 7) **LTY**: 1919 年至 1925 年的长期政府债券利率，来自于 NBER 宏观历史数据库中的美国长期美国债券收益率。1926 年至 2002 年的收益率来自 Ibbotson 的股票、债券、票据和通货膨胀年鉴。2003 年的收益率是长期国库券的均值（25 年及以上）。
- 8) **BOND**: LTY 减去 LTY 过去 12 个月的均值。

- 9) **TERM**: LTY 利率和 TBL 利率的差值。
- 10) **CREDIT**: 评级为 AAA 的公司债利率与 LTY 利率之差。
- 11) **DFY**: 违约收益率利差, 评级为 BAA 的公司债利率与评级为 AAA 的公司债利率之差。
- 12) **DFR**: 违约回报差价, 长期公司债利率与长期政府债利率之差。
- 13) **BM**: 账面市值比。
- 14) **Corpr**: 长期公司债利率。
- 15) **ntis**: 净股权扩张, 股票市场市值净增长的十二个月的移动平均值除以上一年末的股票市场总市值。

数据来源参考 Ivo Welch (2008)。

ii. 市场情绪指标 (7 个)

- 1) **SENT**: Baker and Wurgler (2006) 使用 PCA 方法构建的情绪指标。
- 2) **SENT\_ORTH**: Baker and Wurgler (2006) 使用 PCA 方法构建的情绪指标, 且该因子构建的过程中控制了通货膨胀和产出等宏观变量, 因此该因子与宏观经济趋势是正交的。
- 3) **pdnd**: 股息溢价, 定义为股息支付者和非支付者的平均账面市值比之差。股息溢价很好地解释了公司支付股息倾向的主要历史趋势,
- 4) **ripo**: IPO 首日回报。首次公开募股有时会在首个交易日获得可观的回报, 与投资者的热情有着密切关联。

- 5) **nipo**: 首次公开募股量。首次公开募股的潜在需求对投资者情绪极为敏感。反复无常地打开和关闭的首次公开募股的现象被称为“机会之窗”。这种反复无常可以解释为什么 IPO 数量会出现剧烈波动，在某些时期每月发行超过 100 次，而在其他时期则为每月零发行。
- 6) **cefd**: 封闭式基金折扣。封闭式基金是发行固定数量股票，然后在证券交易所交易的投资公司。封闭式基金的“折价”（或偶尔溢价）是指基金实际持有证券的资产净值与基金市场价格之间的差值。包括 Zweig (1973)、Lee、Shleifer 和 Thaler (1991) 以及 Neal 和 Wheatley (1998) 在内的许多作者都认为，如果散户投资者不成比例地持有封闭式基金，则封闭式股票的平均折价基金可能是一个情绪指数，当散户投资者看跌时，折价会增加。
- 7) **S**: 股票换手率。在一个有卖空限制的市场中，非理性的投资者只有在他们乐观的时候才会参与，从而增加流动性，因此高流动性往往暗示着高涨的投资者情绪。

数据来源参考 Baker & Wurgler (2006)。

### iii. 技术指标 (45 个)

- 1) **MA(1,3)** 技术信号: 一个指标变量，如果标准普尔 500 价格指数的 1 个月移动均线大于或等于 (小于) 标准普尔 500 价格指数的 3 个月移动均线，则取值为 1 (0)。

- 2) **MA(1,6)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 6 个月移动平均线, 则取值为 1 (0)。
- 3) **MA(1,9)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 9 个月移动平均线, 则取值为 1 (0)。
- 4) **MA(1,12)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 12 个月移动平均线, 则取值为 1 (0)。
- 5) **MA(1,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。
- 6) **MA(1,36)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 7) **MA(2,3)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 3 个月移动平均线, 则取值为 1 (0)。

- 8) **MA(2,6)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 6 个月移动平均线, 则取值为 1 (0)。
- 9) **MA(2,9)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 9 个月移动平均线, 则取值为 1 (0)。
- 10) **MA(2,12)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 12 个月移动平均线, 则取值为 1 (0)。
- 11) **MA(2,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。
- 12) **MA(2,36)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 2 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 13) **MA(3,6)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 1 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。

- 14) **MA(3,9)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 3 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 9 (0)。
- 15) **MA(3,12)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 3 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 12 个月移动平均线, 则取值为 1 (0)。
- 16) **MA(3,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 3 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。
- 17) **MA(3,36)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 3 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 18) **MA(6,9)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 6 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 9 个月移动平均线, 则取值为 1 (0)。
- 19) **MA(6,12)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 6 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 12 个月移动平均线, 则取值为 1 (0)。

- 20) **MA(6,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 6 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。
- 21) **MA(6,36)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 6 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 22) **MA(9,12)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 9 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 12 个月移动平均线, 则取值为 1 (0)。
- 23) **MA(9,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 9 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。
- 24) **MA(9,36)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 9 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 25) **MA(12,24)** 技术信号: 一个指标变量, 如果标准普尔 500 价格指数的 12 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 24 个月移动平均线, 则取值为 1 (0)。

- 26) **MA(12,36) 技术信号:** 一个指标变量, 如果标准普尔 500 价格指数的 12 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 27) **MA(24,36) 技术信号:** 一个指标变量, 如果标准普尔 500 价格指数的 24 个月移动平均线大于或等于 (小于) 标准普尔 500 价格指数的 36 个月移动平均线, 则取值为 1 (0)。
- 28) **MOM1 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 1 个月前的收盘价, 则取值为 1 (0)。
- 29) **MOM2 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 2 个月前的收盘价, 则取值为 1 (0)。
- 30) **MOM3 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 3 个月前的收盘价, 则取值为 1 (0)。
- 31) **MOM6 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 6 个月前的收盘价, 则取值为 1 (0)。
- 32) **MOM9 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 9 个月前的收盘价, 则取值为 1 (0)。
- 33) **MOM12 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 12 个月前的收盘价, 则取值为 1 (0)。
- 34) **MOM24 动量:** 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 24 个月前的收盘价, 则取值为 1 (0)。

- 35) **MOM36 动量**: 一个指标变量, 如果标准普尔 500 当月收盘价大于或等于 (小于) 标准普尔 500 价格指数的 36 个月前的收盘价, 则取值为 1 (0)。
- 36) **RSI1 相对强弱指标**: 一个指标变量, 等于 1 个月内收盘涨幅之和除以 1 个月内收盘涨跌幅绝对值之和。
- 37) **RSI2 相对强弱指标**: 一个指标变量, 等于 2 个月内收盘涨幅之和除以 2 个月内收盘涨跌幅绝对值之和。
- 38) **RSI3 相对强弱指标**: 一个指标变量, 等于 3 个月内收盘涨幅之和除以 3 个月内收盘涨跌幅绝对值之和。
- 39) **RSI6 相对强弱指标**: 一个指标变量, 等于 6 个月内收盘涨幅之和除以 6 个月内收盘涨跌幅绝对值之和。
- 40) **RSI9 相对强弱指标**: 一个指标变量, 等于 9 个月内收盘涨幅之和除以 9 个月内收盘涨跌幅绝对值之和。
- 41) **RSI12 相对强弱指标**: 一个指标变量, 等于 12 个月内收盘涨幅之和除以 12 个月内收盘涨跌幅绝对值之和。
- 42) **RSI24 相对强弱指标**: 一个指标变量, 等于 24 个月内收盘涨幅之和除以 24 个月内收盘涨跌幅绝对值之和。
- 43) **RSI36 相对强弱指标**: 一个指标变量, 等于 36 个月内收盘涨幅之和除以 36 个月内收盘涨跌幅绝对值之和。

- 44) **EMA3\_9** 指数平均数指标：一个指标变量，如果标准普尔 500 价格指数的 3 个月指数移动平均线的大于或等于（小于）标准普尔 500 价格指数的 9 个月指数移动平均线，则取值为 1（0）。
- 45) **EMA3\_12** 指数平均数指标：一个指标变量，如果标准普尔 500 价格指数的 3 个月指数移动平均线的大于或等于（小于）标准普尔 500 价格指数的 12 个月指数移动平均线，则取值为 1（0）。

数据来源 WRDS 数据库。

iv. 宏观经济指标（134 个）

- 1) **RPI**: 真实个人收入。
- 2) **W875RX1**: 真实个人收入减转账收据。
- 3) **DPCERA3Mo86SBEA**: 实际个人消费支出。
- 4) **CMRMTSPLx**: 实际制造业和贸易业销售。
- 5) **RETAILx**: 零售和食品服务销售。
- 6) **INDPRO**: 工业生产指数。
- 7) **IPFPNSS**: 工业生产指数：最终产品和非工业用品。
- 8) **IPFINAL**: 工业生产指数：最终产品（市场组）。
- 9) **IPCONGD**: 工业生产指数：消费品。
- 10) **IPDCONGD**: 工业生产指数：耐用消费品。
- 11) **IPNCONG**: 工业生产指数：非耐用消费品。
- 12) **IPBUSEQ**: 工业生产指数：商业设备。

- 13) IPMAT: 工业生产指数: 材料。
- 14) IPDMAT: 工业生产指数: 耐用材料。
- 15) IPNMAT: 工业生产指数: 非耐用材料。
- 16) IPMANSICS: 工业生产指数: 制造业 (标准行业分类)。
- 17) IPB51222S: 工业生产指数: 住宅公用事业。
- 18) IPFUELS: 工业生产指数: 燃料。
- 19) CUMFNS: 产能利用率: 制造业。
- 20) HWI: 美国寻求帮助人数。
- 21) HWIURATIO: 需要帮助/失业人数的比例。
- 22) CLF16OV: 平民劳动力。
- 23) CE16OV: 就业平民数。
- 24) UNRATE: 平民失业率。
- 25) UEMPMEAN: 失业平均时长 (周)。
- 26) UEMPLT5: 失业平民数 (小于 5 周)。
- 27) UEMP5TO14: 失业平民数 (5-14 周)。
- 28) UEMP15OV: 失业平民数 (大于 15 周)。
- 29) UEMP15T26: 失业平民数 (15-26 周)。
- 30) UEMP27OV: 失业平民数 (大于 27 周)。
- 31) CLAIMSx: 失业索赔。
- 32) PAYEMS: 非农就业者。

- 33) USGOOD: 全体员工: 商品生产行业。
- 34) CES1021000001: 全体员工: 采矿业和伐木业: 采矿业。
- 35) USCONS: 全体员工: 建筑业。
- 36) MANEMP: 全体员工: 制造业。
- 37) DMANEMP: 全体员工: 耐用品。
- 38) NDMANEMP: 全体员工: 非耐用品。
- 39) SRVPRD: 全体员工: 服务业。
- 40) USTPU: 全体员工: 贸易、运输和公用事业。
- 41) USWTRADE: 全体员工: 批发。
- 42) USTRAD: 全体员工: 零售。
- 43) USFIRE: 全体员工: 金融活动。
- 44) USGOVT: 全体员工: 政府。
- 45) CES0600000007: 平均每周工作小时数: 商品生产。
- 46) AWOTMAN: 平均每周超时工作小时数: 制造业。
- 47) AWHMAN: 平均每周正常工作小时数: 制造业。
- 48) HOUST: 房屋开工: 新增私有住房总数。
- 49) HOUSTNE: 房屋开工: 东北部。
- 50) HOUSTMW: 房屋开工: 中西部。
- 51) HOUSTS: 房屋开工: 南部。
- 52) HOUSTW: 房屋开工: 西部。

- 53) PERMIT: 新私人房屋许可证。
- 54) PERMITNE: 新私人房屋许可证, 东部。
- 55) PERMITMW: 新私人房屋许可证, 中西部。
- 56) PERMITS: 新私人房屋许可证, 南部。
- 57) PERMITW: 新私人房屋许可证, 西部。
- 58) ACOGNO: 消费品新订单。
- 59) AMDMNOx: 耐用品新订单。
- 60) ANDENOX: 非国防资本货物的新订单。
- 61) AMDMUOX: 耐用品的未成交订单。
- 62) BUSINVx: 企业总库存。
- 63) ISRATIOx: 企业总库存/销售比例。
- 64) M1SL: M1 货币存量。
- 65) M2SL: M2 货币存量。
- 66) M2REAL: 真实 M1 货币存量。
- 67) BOGMBASE: 货币基础。
- 68) TOTRESNS: 存款机构总准备金。
- 69) NONBORRES: 存款机构准备金。
- 70) BUSLOANS: 商业和工业贷款。
- 71) REALLN: 所有商业银行的房地产贷款。
- 72) NONREVSL: 非循环信贷总额。

- 73) CONSPI: 非循环消费信贷与个人收入。
- 74) S&P 500: 标普 500 指数。
- 75) S&P: indust: 标普 500 指数: 工业。
- 76) S&P div yield: 标准普尔综合普通股: 股息收益。
- 77) S&P PE ratio: 标准普尔综合普通股: 市盈率。
- 78) FEDFUNDS: 有效联邦基金利率。
- 79) CP3Mx: 3 个月 AA 金融商业票据利率。
- 80) TB3MS: 3 个月国库券利率。
- 81) TB6MS: 6 个月国库券利率。
- 82) GS1: 1 年国库券利率。
- 83) GS5: 5 年国库券利率。
- 84) GS10: 10 年国库券利率。
- 85) AAA: 穆迪 Aaa 公司债券收益率。
- 86) BAA: 穆迪 Baa 公司债券收益率。
- 87) COMPAPFFx: 3 个月期商业票据减去联邦基金。
- 88) TB3SMFFM: 3 个月国债 C 减去联邦基金。
- 89) TB6SMFFM: 6 个月国债 C 减去联邦基金。
- 90) T1YFFM: 1 年国债 C 减去联邦基金。
- 91) T5YFFM: 5 年国债 C 减去联邦基金。
- 92) T10YFFM: 10 年国债 C 减去联邦基金。

- 93) AAFFM: 穆迪 Aaa 公司债券收益率减去联邦基金。
- 94) BAAFFM: 穆迪 Baa 公司债券收益率减去联邦基金。
- 95) TWEXAFEGSMTHx: 贸易加权美元指数。
- 96) EXSZUSx: 瑞士 / 美国外汇汇率。
- 97) EXJPUSx: 日本 / 美国外汇汇率。
- 98) EXUSUKx: 美国 / 英国外汇汇率。
- 99) EXCAUSx: 加拿大 / 美国外汇汇率。
- 100) WPSFD49207: 生产者价格指数 (PPI): 制成品。
- 101) WPSFD49502: 生产者价格指数 (PPI): 制成消费品。
- 102) WPSID61: 生产者价格指数 (PPI): 中间产品。
- 103) WPSID62: 生产者价格指数 (PPI): 原材料。
- 104) OILPRICEx: 生产者价格指数 (PPI): 原油、WTI 和库欣。
- 105) PPICMM: 生产者价格指数 (PPI): 金属和金属制品。
- 106) CPIAUCSL: 居民消费价格指数: 所有项目。
- 107) CPIAPPSL: 居民消费价格指数: 服装。
- 108) CPITRNSL: 居民消费价格指数: 交通。
- 109) CPIMEDSL: 居民消费价格指数: 医药。
- 110) CUSR0000SAC: 居民消费价格指数: 日用品。
- 111) CUSR0000SAD: 居民消费价格指数: 耐用品。
- 112) CUSR0000SAS: 居民消费价格指数: 服务。

- 113) CPIULFSL: 居民消费价格指数: 所有项目减去食物。
- 114) CUSR0000SAoL2: 居民消费价格指数: 所有项目减去住。
- 115) CUSR0000SAoL5: 居民消费价格指数: 所有项目减去医药。
- 116) PCEPI: 个人开支: 连锁指数。
- 117) DDURRG3Mo86SBEA: 个人开支: 耐用品。
- 118) DNDGRG3Mo86SBEA: 个人开支: 非耐用品。
- 119) DSERRG3Mo86SBEA: 个人开支: 服务。
- 120) CESo6000000008: 平均每小时收入: 商品生产。
- 121) CES20000000008: 平均每小时收入: 建筑业。
- 122) CES30000000008: 平均每小时收入: 制造业。
- 123) UMCSENTx: 消费者信心指数。
- 124) DTCOLNVHFNm: 未偿还消费机动车贷款。
- 125) DTCTHFNm: 未偿消费贷款和租赁总额。
- 126) INVEST: 所有商业银行的银行信贷证券。
- 127) VIXCLSx: 恐慌指数。
- 128) Fred Factor1: 通过以上 127 个因子运算得到。
- 129) Fred Factor2: 通过以上 127 个因子运算得到。
- 130) Fred Factor3: 通过以上 127 个因子运算得到。
- 131) Fred Factor4: 通过以上 127 个因子运算得到。
- 132) Fred Factor5: 通过以上 127 个因子运算得到。

133) Fred Factor6: 通过以上 127 个因子运算得到。

134) Fred Factor7: 通过以上 127 个因子运算得到。

数据来源参考 McCracken (2016)。

v. 另一来源的宏观因子 (9 个)

1) F1: 对美国长期债券收益率分解得到的分量。

2) F2: 对美国长期债券收益率分解得到的分量。

3) F3: 对美国长期债券收益率分解得到的分量。

4) F4: 对美国长期债券收益率分解得到的分量。

5) F5: 对美国长期债券收益率分解得到的分量。

6) F6: 对美国长期债券收益率分解得到的分量。

7) F7: 对美国长期债券收益率分解得到的分量。

8) F8: 对美国长期债券收益率分解得到的分量。

9)  $F1^3$ : F1 因子的三次方。

数据来源参考 Ludvigson (2009)。

vi. 中期选举指标 (11 个)

1) mid\_-12: 该月及之前的 12 个月是否有中期选举。

2) mid\_-6: 该月及之前的 6 个月是否有中期选举。

3) mid\_-3: 该月及之前的 3 个月是否有中期选举。

4) mid\_-2: 该月及之前的 2 个月是否有中期选举。

5) mid\_-1: 该月及之前的 1 个月是否有中期选举。

- 6) mid\_0: 该月是否有中期选举。
- 7) mid\_1: 该月及之后的 1 个月是否有中期选举。
- 8) mid\_2: 该月及之后的 2 个月是否有中期选举。
- 9) mid\_3: 该月及之后的 3 个月是否有中期选举。
- 10) mid\_6: 该月及之后的 6 个月是否有中期选举。
- 11) mid\_12: 该月及之后的 12 个月是否有中期选举。

以上变量根据中期选举所在月份计算得到。

#### 4.1.2 自变量描述性统计

基于上述变量，对每个变量分别计算均值、标准差、偏度、峰度、自回归系数、前 10%分位数、中位数、后 10%分位数。各统计量的定义如下所示：

##### 1) 均值 (Mean):

本文中使用的均值为算数平均值，也就是将数据集中目标特征的所有数值相加并除以数据集的大小。均值是一种用于描述数据集中的集中趋势的统计量，是统计学中最基本和常用的测量中心的指标之一，它可以反映数据集中的典型数值水平，也可以用来比较不同数据集之间的差异。

计算样本数据的算术平均值，具体计算公式如下：

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{\sum_1^n x_i}{n}$$

##### 2) 标准差 (Standard deviation):

此处指样本标准差。为样本总体各单位标准值与其均值离差平方的算术平均值的平方根用于描述一组数据的离散程度。它是数据与其平均值之间的差异的平方的平均值的平方

根。标准差越大，表示数据越分散，反之则表示数据越聚集在平均值附近。标准差的计算对于理解数据的分布及其可靠性具有重要意义，因为它可以帮助识别数据中的异常值或离群值，并判断数据是否具有一定的统计意义。

具体计算公式如下：

$$S = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

其中 $x_i$ 为第*i*个样本的值， $\bar{x}$ 为样本均值， $n$ 为样本数。

### 3) 偏度 (Skew):

偏度是样本的三阶标准化矩，用于描述概率分布的偏态或偏斜程度。它是概率密度函数的三阶中心矩与标准差的比值，反映了数据分布相对于平均值的不对称程度。当偏度为0时，表示数据分布是对称的，正偏斜表示数据分布偏向于左侧（即均值左侧的尾部更长），而负偏斜表示数据分布偏向于右侧。偏度的计算对于理解数据的分布形状和特征具有重要意义，因为它可以帮助比较不同数据集的偏态差异，并识别可能存在的异常值。

$$S_k = \frac{\mu_3}{\frac{\mu_2^{\frac{3}{2}}}{\sigma^3}}$$

其中 $\mu_3$ 为三阶中心矩， $\sigma$ 是标准差。

### 4) 峰度 (Kurt):

峰度用于描述概率分布的峰态或尖峰程度。它是概率密度函数在其均值处峰值的相对尖度的测量，也就是随机变量的概率密度函数的四阶中心矩与方差平方的比值。峰度值大于3表示分布的峰比正态分布更尖，而峰度值小于3则表示分布的峰比正态分布更平坦。峰度的计算对于理解数据的分布形状和特征具有重要意义，因为它可以帮助比较不同数据

集的峰态差异，判断数据的分布是否对称或偏态，以及识别可能存在的异常值。

具体计算公式如下：

$$Kurtosis = \frac{\mu_4}{\sigma^4}$$

其中  $\mu_4$  为四阶中心距， $\sigma$  是标准差。

#### 5) 一阶自回归模型的自回归系数 (AR(1)):

一阶自回归模型 (AR(1)模型) 是时间序列分析中常用的模型之一，它是指当前时间点的值与前一时间点的值之间存在线性关系，自回归系数表示前一时间点的值对当前时间点的值的影响程度，即一个单位的变化对当前值的影响程度。当自回归系数大于 0 时，表示前一时间点的值越大，当前时间点的值也越大，反之亦然。当自回归系数等于 0 时，表示当前时间点的值不受前一时间点的值的影响，即不存在自回归效应。自回归系数的估计可以通过最大似然估计或最小二乘法等方法获得，它的大小和符号对于时间序列分析和预测具有重要的意义。

AR(1) 模型的表达式为  $X_t = aX_{t-1} + \epsilon_t$ ，其中  $\{\epsilon_t\} \sim N(0, \sigma^2)$ 。其中  $X_t$  表示当前时间点的值， $X_{t-1}$  表示前一时间点的值， $\epsilon_t$  表示误差项。

#### 6) 前 10%分位数 (top 10%)

前 10%分位数，也称为第一分位数，是指将一组数据按照从小到大的顺序排列后，将其分为 10 份，而第一份中最大的数值即为前 10%分位数。前 10%分位数在数据分析中非常有用，因为它可以帮助了解数据集的中心趋势和离散程度。通过计算前 10%分位数，可以得到一组数据中最小的 10%的值，这有助于了解这个数据集的最小值的情况。同时，还可以结合其他的分位数，比如中位数和第三分位数，来了解数据的整体分布情况。

## 7) 中位数 (Median)

中位数是一组数据中的一个特殊值，它将这组数据分为相等的两部分。具体来说，中位数是将一组数据按照从小到大的顺序排列后，位于中间的那个数，如果这组数据的数量是偶数，那么中位数就是中间两个数的平均数。中位数是一种描述数据集中心位置的指标，相较于均值，中位数更能反映数据的实际情况。因为中位数不受数据的分布情况影响，可以避免极端值的干扰。在数据分析中，中位数可以用来描述数据的集中趋势和离散程度，以及数据的异常情况。

## 8) 后 10%分位数 (bottom 10%)

后 10%分位数，也称为第九分位数，是指将一组数据按照从小到大的顺序排列后，将其分为 10 份，而第九份的最大数值即为后 10%分位数。后 10%分位数在数据分析中用于描述数据集中最大的 10%数值的情况。它是数据分布的重要指标之一，与前 10%分位数一起可以用来判断数据的离散程度。如果数据的后 10%分位数与前 10%分位数相差很大，那么这个数据集就具有很大的离散程度。而如果两者相差不大，那么这个数据集就具有相对较小的离散程度。

标准差是一种反映数据集中数据分散程度的统计量。当一组数据的标准差值计算得到较小时，通常会出现较多 0.00 的取值，这是因为四舍五入后取整导致的。这种情况通常发生在数据集中的差异性较小的情况下。除了标准差，本文中其他一些统计量，比如均值和中位数，也可能出现类似的情况。这是由于这些统计量是基于原始数据计算得出的，当数据集合的差异性较小，这些统计量的取值也会比较接近，导致四舍五入后出现类似的取值。

自变量的具体统计结果见表 1 所示。

表 1 自变量统计结果

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
DP	-3.65	0.42	0.08	-0.99	1.00	-3.04	-3.61	-4.15
DY	-3.64	0.42	0.07	-0.97	1.00	-3.04	-3.60	-4.14
EP	-2.87	0.45	-0.52	2.04	0.99	-2.21	-2.89	-3.36
VOL	14.66	5.29	0.75	0.38	0.96	21.62	13.69	8.55
BILL	-0.04	0.93	-0.11	3.38	0.90	0.93	-0.01	-1.15
TBL	4.54	3.39	0.60	0.36	0.99	8.59	4.89	0.09
BOND	-0.03	0.56	-0.12	2.02	0.86	0.61	-0.03	-0.67
LTY	6.43	2.94	0.38	-0.13	1.00	10.57	6.25	2.60
TERM	1.89	1.47	-0.42	-0.07	0.95	3.81	1.96	-0.01
CREDIT	0.83	0.42	0.31	-0.09	0.81	1.36	0.76	0.34
DFY	1.07	0.43	1.89	4.76	0.96	1.63	0.94	0.67
DFR	0.02	1.56	-0.67	7.24	-0.04	1.64	0.06	-1.51
BM	47.43	27.34	0.88	-0.37	1.00	92.75	35.46	19.57
corpr	0.69	2.74	0.40	3.33	0.09	3.61	0.66	-2.40
ntis	0.83	2.06	-0.43	-0.24	0.98	3.40	1.09	-2.14
SENT	0.04	0.96	0.44	0.91	0.99	1.22	-0.08	-1.10
SENT_ORTH	0.04	0.97	0.17	1.27	0.98	1.05	0.02	-1.16
pdnd	-5.91	12.85	0.32	0.66	0.97	14.05	-7.19	-20.71
ripo	17.85	19.18	2.12	6.31	0.61	37.80	13.50	0.71
nipo	27.56	24.57	1.37	2.07	0.85	62.00	20.00	3.00
cefd	8.72	6.73	-0.01	0.03	0.96	18.21	9.11	0.58
s	0.17	0.09	0.99	0.41	0.99	0.30	0.15	0.08
MA1_3	0.63	0.48	-0.53	-1.72	0.23	1.00	1.00	0.00
MA1_6	0.67	0.47	-0.74	-1.46	0.54	1.00	1.00	0.00
MA1_9	0.70	0.46	-0.85	-1.27	0.69	1.00	1.00	0.00
MA1_12	0.73	0.45	-1.01	-0.98	0.77	1.00	1.00	0.00
MA1_24	0.77	0.42	-1.32	-0.26	0.78	1.00	1.00	0.00
MA1_36	0.81	0.39	-1.57	0.48	0.84	1.00	1.00	0.00
MA2_3	0.62	0.49	-0.49	-1.76	0.27	1.00	1.00	0.00
MA2_6	0.69	0.46	-0.83	-1.31	0.64	1.00	1.00	0.00

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
MA2_9	0.71	0.46	-0.90	-1.19	0.75	1.00	1.00	0.00
MA2_12	0.72	0.45	-1.00	-1.00	0.81	1.00	1.00	0.00
MA2_24	0.77	0.42	-1.31	-0.29	0.81	1.00	1.00	0.00
MA2_36	0.81	0.39	-1.57	0.48	0.90	1.00	1.00	0.00
MA3_6	0.68	0.47	-0.77	-1.41	0.66	1.00	1.00	0.00
MA3_9	0.71	0.46	-0.90	-1.19	0.78	1.00	1.00	0.00
MA3_12	0.72	0.45	-1.00	-1.00	0.81	1.00	1.00	0.00
MA3_24	0.78	0.42	-1.33	-0.24	0.88	1.00	1.00	0.00
MA3_36	0.81	0.39	-1.56	0.44	0.92	1.00	1.00	0.00
MA6_9	0.70	0.46	-0.85	-1.27	0.79	1.00	1.00	0.00
MA6_12	0.72	0.45	-0.97	-1.06	0.85	1.00	1.00	0.00
MA6_24	0.78	0.41	-1.38	-0.09	0.91	1.00	1.00	0.00
MA6_36	0.80	0.40	-1.52	0.32	0.93	1.00	1.00	0.00
MA9_12	0.72	0.45	-0.97	-1.06	0.84	1.00	1.00	0.00
MA9_24	0.79	0.41	-1.42	0.01	0.91	1.00	1.00	0.00
MA9_36	0.81	0.39	-1.57	0.48	0.93	1.00	1.00	0.00
MA12_24	0.79	0.41	-1.41	-0.02	0.92	1.00	1.00	0.00
MA12_36	0.81	0.39	-1.60	0.56	0.94	1.00	1.00	0.00
MA24_36	0.80	0.40	-1.50	0.25	0.92	1.00	1.00	0.00
MOM1	0.60	0.49	-0.42	-1.83	0.05	1.00	1.00	0.00
MOM2	0.63	0.48	-0.54	-1.71	0.31	1.00	1.00	0.00
MOM3	0.65	0.48	-0.64	-1.59	0.49	1.00	1.00	0.00
MOM6	0.69	0.46	-0.81	-1.34	0.68	1.00	1.00	0.00
MOM9	0.72	0.45	-1.00	-1.00	0.75	1.00	1.00	0.00
MOM12	0.75	0.43	-1.15	-0.68	0.80	1.00	1.00	0.00
MOM24	0.80	0.40	-1.51	0.28	0.84	1.00	1.00	0.00
MOM36	0.81	0.39	-1.56	0.44	0.86	1.00	1.00	0.00
RSI1	51.54	50.01	-0.06	-2.00	-0.33	100.00	100.00	0.00
RSI2	51.79	34.58	-0.04	-1.21	0.01	100.00	51.22	0.00
RSI3	51.32	23.83	0.06	-0.22	0.06	83.12	50.00	21.96
RSI6	50.74	11.67	0.30	0.58	0.02	65.26	50.21	36.59

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
RSI9	50.33	7.57	0.62	1.93	-0.01	59.87	49.71	41.67
RSI12	50.19	5.60	0.18	0.73	-0.00	56.94	49.99	43.57
RSI24	50.03	2.82	0.22	0.67	0.03	53.45	49.92	46.67
RSI36	50.01	1.93	0.26	1.03	0.02	52.36	49.92	47.89
EMA3_9	0.52	0.50	-0.07	-2.00	0.17	1.00	1.00	0.00
EMA3_12	0.52	0.50	-0.10	-2.00	0.20	1.00	1.00	0.00
Macro1	0.00	0.01	-0.19	14.53	-0.06	0.01	0.00	-0.00
Macro2	0.00	0.01	-0.85	20.32	0.06	0.01	0.00	-0.00
Macro3	0.00	0.01	-0.16	25.75	-0.05	0.01	0.00	-0.00
Macro4	0.00	0.01	-3.06	51.75	-0.05	0.01	0.00	-0.01
Macro5	0.00	0.01	-1.16	24.44	-0.05	0.02	0.00	-0.01
Macro6	0.00	0.01	-6.31	100.46	0.26	0.01	0.00	-0.01
Macro7	0.00	0.01	-5.39	89.30	0.18	0.01	0.00	-0.01
Macro8	0.00	0.01	-4.86	84.38	0.12	0.01	0.00	-0.01
Macro9	0.00	0.01	-4.11	71.18	0.07	0.01	0.00	-0.01
Macro10	0.00	0.03	-3.96	74.54	0.11	0.02	0.00	-0.02
Macro11	0.00	0.01	-0.41	2.14	-0.19	0.01	0.00	-0.01
Macro12	0.00	0.02	-3.80	58.41	0.20	0.02	0.00	-0.01
Macro13	0.00	0.01	-5.61	77.57	0.27	0.01	0.00	-0.01
Macro14	0.00	0.02	-5.54	78.81	0.34	0.02	0.00	-0.01
Macro15	0.00	0.01	-2.02	15.53	0.08	0.01	0.00	-0.01
Macro16	0.00	0.01	-6.38	104.54	0.23	0.01	0.00	-0.01
Macro17	0.00	0.04	-0.01	1.46	-0.22	0.05	0.00	-0.05
Macro18	0.00	0.02	-0.38	13.45	-0.07	0.02	0.00	-0.02
Macro19	-0.02	0.88	-6.29	101.02	0.22	0.71	0.04	-0.65
Macro20	12.24	195.53	-0.15	5.38	-0.12	224.30	10.00	-201.30
Macro21	0.00	0.04	-0.31	4.73	0.13	0.04	0.00	-0.04
Macro22	0.00	0.00	-0.13	4.37	-0.13	0.00	0.00	-0.00
Macro23	0.00	0.00	-0.64	24.65	0.11	0.00	0.00	-0.00
Macro24	-0.01	0.22	-0.12	18.94	0.21	0.20	0.00	-0.20
Macro25	0.04	0.75	-0.16	6.00	0.06	0.80	0.00	-0.70

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
Macro26	0.00	0.06	0.42	6.95	-0.34	0.07	0.00	-0.07
Macro27	-0.00	0.06	-0.48	9.35	-0.17	0.07	-0.00	-0.07
Macro28	0.00	0.05	0.95	5.61	0.25	0.06	-0.00	-0.06
Macro29	0.00	0.08	0.01	10.80	0.08	0.09	0.00	-0.09
Macro30	0.00	0.07	0.76	3.55	0.11	0.09	-0.00	-0.08
Macro31	-0.00	0.06	0.91	21.92	0.14	0.06	-0.00	-0.06
Macro32	0.00	0.00	-3.02	48.92	0.39	0.00	0.00	-0.00
Macro33	-0.00	0.01	-3.40	47.06	0.36	0.00	0.00	-0.01
Macro34	-0.00	0.01	-1.01	13.88	0.33	0.01	0.00	-0.01
Macro35	0.00	0.01	-0.34	18.32	0.13	0.01	0.00	-0.01
Macro36	-0.00	0.01	-3.93	47.32	0.43	0.00	0.00	-0.01
Macro37	-0.00	0.01	-3.71	42.73	0.40	0.00	0.00	-0.01
Macro38	-0.00	0.00	-1.93	24.23	0.42	0.00	-0.00	-0.00
Macro39	0.00	0.00	-1.13	28.27	0.37	0.00	0.00	-0.00
Macro40	0.00	0.00	-1.33	18.11	0.43	0.00	0.00	-0.00
Macro41	0.00	0.00	-0.87	3.41	0.60	0.00	0.00	-0.00
Macro42	0.00	0.00	-0.54	22.60	0.31	0.00	0.00	-0.00
Macro43	0.00	0.00	-0.53	1.37	0.74	0.00	0.00	-0.00
Macro44	0.00	0.00	0.20	17.24	0.08	0.00	0.00	-0.00
Macro45	40.31	0.66	-0.29	0.17	0.91	41.20	40.30	39.50
Macro46	0.00	0.14	-0.69	10.27	-0.25	0.10	0.00	-0.10
Macro47	40.81	0.76	-0.45	0.23	0.93	41.80	40.80	39.80
Macro48	7.22	0.33	-0.98	0.92	0.97	7.59	7.29	6.81
Macro49	5.00	0.40	-0.26	0.02	0.85	5.54	5.00	4.47
Macro50	5.52	0.43	-0.75	0.04	0.91	5.99	5.62	4.91
Macro51	6.44	0.31	-0.69	0.39	0.94	6.80	6.47	6.01
Macro52	5.79	0.38	-1.04	0.96	0.94	6.22	5.86	5.29
Macro53	7.19	0.32	-0.85	0.50	0.98	7.55	7.24	6.76
Macro54	5.03	0.38	-0.13	-0.03	0.92	5.56	5.03	4.57
Macro55	5.49	0.40	-0.68	-0.19	0.96	5.93	5.56	4.92
Macro56	6.37	0.33	-0.50	-0.12	0.98	6.78	6.40	5.91

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
Macro57	5.82	0.38	-1.01	0.97	0.97	6.25	5.88	5.33
Macro58	0.01	0.02	-1.04	21.68	0.31	0.03	0.01	-0.01
Macro59	0.00	0.04	-0.54	5.54	-0.25	0.05	0.00	-0.04
Macro60	0.00	0.09	-0.10	5.14	-0.39	0.10	0.00	-0.10
Macro61	0.00	0.01	0.66	1.72	0.62	0.02	0.00	-0.01
Macro62	0.00	0.01	-0.52	3.76	0.69	0.01	0.00	-0.00
Macro63	-0.00	0.02	1.93	39.98	-0.02	0.02	0.00	-0.02
Macro64	0.00	0.01	0.21	10.83	-0.38	0.01	0.00	-0.01
Macro65	0.00	0.00	-0.07	10.16	-0.27	0.00	0.00	-0.00
Macro66	0.00	0.01	1.92	12.77	0.50	0.01	0.00	-0.00
Macro67	0.00	0.02	0.77	13.33	-0.02	0.02	0.00	-0.02
Macro68	0.00	0.05	0.10	13.02	-0.10	0.05	0.00	-0.06
Macro69	0.00	0.05	-0.12	10.09	-0.18	0.05	0.00	-0.05
Macro70	0.00	0.01	-0.09	5.61	-0.26	0.01	0.00	-0.01
Macro71	0.00	0.00	0.53	11.73	-0.25	0.00	0.00	-0.00
Macro72	-0.00	0.00	-0.60	12.02	-0.44	0.00	-0.00	-0.00
Macro73	0.00	0.00	1.24	18.34	0.10	0.00	0.00	-0.00
Macro74	0.01	0.04	-1.24	5.28	0.23	0.04	0.01	-0.03
Macro75	0.01	0.04	-1.17	4.75	0.24	0.05	0.01	-0.03
Macro76	-0.00	0.12	0.66	6.03	0.27	0.12	-0.01	-0.11
Macro77	0.00	0.05	-0.05	5.38	0.51	0.05	0.00	-0.05
Macro78	-0.01	0.37	-0.38	6.62	0.51	0.33	0.01	-0.38
Macro79	-0.01	0.39	-0.71	5.96	0.42	0.36	0.00	-0.38
Macro80	-0.00	0.36	-0.35	6.36	0.34	0.31	0.00	-0.37
Macro81	-0.00	0.37	-0.36	7.53	0.35	0.32	0.00	-0.40
Macro82	-0.00	0.40	-0.11	7.40	0.35	0.39	0.00	-0.41
Macro83	-0.01	0.34	-0.39	6.30	0.34	0.35	-0.01	-0.36
Macro84	-0.01	0.29	-0.40	5.51	0.30	0.33	-0.01	-0.31
Macro85	-0.00	0.23	-0.23	4.95	0.31	0.24	-0.01	-0.22
Macro86	-0.00	0.23	0.60	6.15	0.39	0.24	-0.01	-0.24
Macro87	0.04	0.40	-1.30	9.68	0.75	0.43	0.06	-0.38

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
Macro88	-0.51	0.72	-2.54	9.04	0.88	0.01	-0.27	-1.27
Macro89	-0.38	0.77	-2.60	9.99	0.89	0.23	-0.14	-1.22
Macro90	-0.00	0.78	-2.22	8.73	0.86	0.67	0.11	-0.83
Macro91	0.72	1.41	-1.37	3.17	0.94	2.22	0.92	-0.89
Macro92	1.09	1.69	-1.19	2.18	0.95	3.00	1.31	-0.84
Macro93	2.21	1.99	-1.05	1.61	0.97	4.55	2.42	-0.20
Macro94	3.27	2.05	-0.71	0.83	0.97	5.51	3.48	0.78
Macro95	-0.00	0.02	-0.02	0.79	0.34	0.02	-0.00	-0.02
Macro96	-0.00	0.03	-0.10	1.32	0.27	0.03	-0.00	-0.04
Macro97	-0.00	0.03	-0.48	1.54	0.32	0.03	-0.00	-0.03
Macro98	-0.00	0.02	-0.44	2.40	0.33	0.02	0.00	-0.03
Macro99	0.00	0.01	0.51	7.00	0.27	0.02	-0.00	-0.02
Macro100	0.00	0.01	0.33	6.52	-0.42	0.01	-0.00	-0.01
Macro101	0.00	0.01	0.46	7.99	-0.42	0.01	-0.00	-0.01
Macro102	0.00	0.01	-0.83	9.98	-0.36	0.01	-0.00	-0.01
Macro103	0.00	0.04	-0.14	10.29	-0.43	0.04	-0.00	-0.04
Macro104	-0.00	0.09	0.23	2.42	-0.30	0.11	0.00	-0.10
Macro105	0.00	0.03	-0.31	3.64	-0.32	0.04	-0.00	-0.04
Macro106	0.00	0.00	0.05	5.08	-0.25	0.00	-0.00	-0.00
Macro107	0.00	0.01	0.54	5.47	-0.36	0.01	-0.00	-0.01
Macro108	0.00	0.01	-0.22	7.01	-0.12	0.01	-0.00	-0.01
Macro109	-0.00	0.00	-0.30	8.11	-0.53	0.00	-0.00	-0.00
Macro110	0.00	0.01	-0.34	6.16	-0.21	0.01	-0.00	-0.01
Macro111	0.00	0.00	0.33	16.59	-0.20	0.00	-0.00	-0.00
Macro112	0.00	0.00	0.49	4.97	-0.47	0.00	-0.00	-0.00
Macro113	0.00	0.00	-0.24	2.71	-0.21	0.00	-0.00	-0.00
Macro114	0.00	0.00	-0.22	5.08	-0.24	0.00	-0.00	-0.00
Macro115	0.00	0.00	-0.08	3.34	-0.24	0.00	-0.00	-0.00
Macro116	0.00	0.00	-0.12	1.91	-0.26	0.00	0.00	-0.00
Macro117	0.00	0.00	0.15	1.47	-0.36	0.00	-0.00	-0.00
Macro118	0.00	0.01	-0.41	4.93	-0.20	0.01	0.00	-0.01

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
Macro119	-0.00	0.00	1.16	14.01	-0.47	0.00	-0.00	-0.00
Macro120	-0.00	0.00	-0.17	7.59	-0.60	0.00	-0.00	-0.00
Macro121	0.00	0.01	-0.25	6.59	-0.64	0.01	-0.00	-0.01
Macro122	-0.00	0.00	0.28	6.06	-0.55	0.00	-0.00	-0.00
Macro123	-0.16	3.73	-0.21	1.96	0.00	4.20	-0.20	-4.35
Macro124	0.00	0.02	0.36	8.27	-0.38	0.02	0.00	-0.02
Macro125	-0.00	0.01	-0.04	12.44	-0.40	0.01	0.00	-0.01
Macro126	-0.00	0.01	-0.54	6.70	-0.33	0.01	0.00	-0.01
Macro127	19.79	7.21	1.93	6.19	0.83	28.36	18.07	12.84
Fred Factor1	0.02	0.43	4.64	66.72	0.44	0.35	-0.01	-0.33
Fred Factor2	-0.01	0.28	-0.72	5.91	0.19	0.27	-0.00	-0.31
Fred Factor3	0.00	0.28	0.10	2.13	0.59	0.33	0.01	-0.34
Fred Factor4	0.00	0.25	-0.13	2.83	0.47	0.28	0.01	-0.27
Fred Factor5	-0.01	0.23	-0.94	7.41	0.69	0.23	0.03	-0.25
Fred Factor6	0.00	0.19	-0.06	3.37	0.34	0.23	-0.01	-0.19
Fred Factor7	-0.00	0.18	0.86	4.15	0.43	0.21	-0.01	-0.21
F1	0.02	0.40	1.32	4.32	0.73	0.43	-0.02	-0.39
F2	0.00	0.28	0.19	3.73	0.45	0.28	0.01	-0.33
F3	-0.00	0.27	-0.06	3.80	0.26	0.29	0.01	-0.33
F4	-0.00	0.24	-0.00	1.60	0.55	0.28	-0.01	-0.27
F5	0.00	0.22	1.17	4.37	0.55	0.25	-0.02	-0.22
F6	0.00	0.20	-0.01	0.82	0.62	0.25	0.00	-0.24
F7	0.00	0.18	-1.15	5.74	0.18	0.20	0.01	-0.19
F8	0.01	0.15	-0.12	1.07	0.26	0.21	0.01	-0.17
F1^3	0.09	0.71	9.87	130.98	0.45	0.08	-0.00	-0.06
mid_-12	0.24	0.43	1.22	-0.52	0.89	1.00	0.00	0.00
mid_-6	0.12	0.33	2.34	3.48	0.81	1.00	0.00	0.00
mid_-3	0.06	0.24	3.71	11.78	0.65	0.00	0.00	0.00
mid_-2	0.04	0.20	4.70	20.13	0.48	0.00	0.00	0.00
mid_-1	0.02	0.14	6.86	45.22	-0.02	0.00	0.00	0.00
mid_0	0.02	0.14	6.86	45.22	-0.02	0.00	0.00	0.00

	Mean	Standard deviation	Skew	Kurt	AR(1)	top 10%	Median	Bottom 10%
mid_1	0.02	0.14	6.86	45.22	-0.02	0.00	0.00	0.00
mid_2	0.04	0.20	4.70	20.13	0.48	0.00	0.00	0.00
mid_3	0.06	0.24	3.71	11.78	0.65	0.00	0.00	0.00
mid_6	0.12	0.33	2.34	3.48	0.81	1.00	0.00	0.00
mid_12	0.24	0.43	1.22	-0.52	0.89	1.00	0.00	0.00

#### 4.2 检验指标

Goyal 和 Welch (2008) 指出, 良好的样本内结果并不一定能保证准确的样本外预测, 因为可能存在过度拟合的问题。在机器学习中, 过度拟合通常指模型在训练数据上表现良好, 但在测试集上表现差的情况。过度拟合通常发生在模型过于复杂或数据量不足的情况下。因此, 样本内表现并不能充分代表模型的准确性, 为准确评估模型性能, 评估预测变量的样本外预测准确性至关重要。

在投资实践中, 样本外表现对于投资决策至关重要。投资者希望了解他们投资的资产未来的表现, 而不仅仅是过去的表现。具有良好样本外表现的模型可以帮助投资者做出更准确的投资决策, 并减少投资风险。因此, 从投资实践的角度考虑, 评估预测变量的样本外预测准确性更值得关注, 因为它能够提供更准确的预测结果。

样本外预测的规则是, 只有截至时间  $t$  之前可用的信息可以用于预测  $t+1$  的股票回报。  $t+1$  的预测由  $\hat{R}_{t+1} = \hat{\alpha}_t + \hat{\beta}_t x_t$  给出。其中  $\hat{\alpha}_t$  和  $\hat{\beta}_t$  通过基于预测变量  $\{x_j\}_{j=1}^{t-1}$  和  $\{R_j\}_{j=2}^t$  来估计。根据这样的规则, 采用递归估计窗口在样本外预测中, 即每次完成预测时, 最新的观测值将被放入估计样本中, 并且将重新估计模型参数以进行下一次预测。

$R_{OS}^2$  统计量在本文中用作样本外预测性能的评估工具。此外, 本文还使用流行的 Clark 和 West(2007) MSFE 调整统计量来检验可预测性是否显著。

#### 4.2.1 样本外 R 方统计量

样本外 R 方  $R_{OS}^2$  是一种用于评估机器学习模型预测能力的指标。与传统的  $R^2$  不同， $R_{OS}^2$  可以评估模型在新数据上的表现，因此在评估模型的泛化能力方面更具有参考价值。

$R_{OS}^2$  统计量的计算公式如下：

$$R_{OS}^2 = 1 - \frac{\sum_{t=m+1}^T \hat{e}_{t|t-1}^2}{\sum_{t=m+1}^T \hat{e}_{t|t-1}^{PM}^2}$$

其中  $m$  是初始维持样本外周期的最后一个观测值， $T$  是可用观测值的总数， $\hat{e}_{t|t-1} = r_t - \hat{r}_{t|t-1}$  通常表示某个机器学习模型的预测误差， $\hat{e}_{t|t-1}^{PM} = r_t - \hat{r}_{t|t-1}^{PM}$  是基准模型—历史均值模型的预测误差。 $R_{OS}^2$  类似于样本内  $R^2$  统计量，衡量某种预测方式相对于历史平均的基准预测方法带来的均方预测误差（MSE）减小的比例。如果  $R_{OS}^2 > 0$ ，则表示预测模型的性能优于历史均值基准模型。

当然，由于月度股票回报率本质上包含有限的可预测成分，因此  $R_{OS}^2$  统计量的值普遍较小。这意味着股票市场的波动性和随机性非常高，很难准确地预测未来的市场表现。尽管如此，Campbell 和 Thompson（2008）认为每月  $R_{OS}^2$  值达到大约 0.5% 则可表明市场超额回报的可预测性在经济上具有显著程度。

#### 4.2.2 CW-t 检验

Clark and West(2007) 的修正 t 统计量是一种经典的统计检验方法，用于比较两个时间序列预测模型的预测精度。该方法在计算两个模型的预测误差时，考虑了误差的自相关性和异方差性，可以提高比较结果的准确性。在 CW-t 检验中，使用  $R_{OS}^2 \leq 0$  作为零假设，使用  $R_{OS}^2 > 0$  的作为备选假设，进行检验。

统计检验量为：

$$\frac{\sqrt{P\bar{f}}}{[\hat{f}_{t+1} - \bar{f} \text{的样本方差}]^{\frac{1}{2}}}$$

其中,  $\hat{f}_{t+1} = \hat{e}_{1,t+1}^2 - \hat{e}_{2,t+1}^2$ ,  $\bar{f} = P^{-1}\Sigma_{t=R}^T \hat{f}_{t+1}$ ,  $P$  为样本外的预测数量。

$\hat{e}_{1,t+1}^2$ 和 $\hat{e}_{2,t+1}^2$ 分别为两个预测模型残差的平方, 也即 $\hat{e}_{1,t+1} = y_{1,t+1} - \hat{y}_{1,t+1}$ ,  $\hat{e}_{2,t+1} = y_{2,t+1} - \hat{y}_{2,t+1}$ ,  $t \in [1, P]$ , 该检验统计量在零假设下服从  $t$  分布。基于该统计检验量对两模型的预测精度是否有显著差异进行检验。

### 4.3 样本外预测结果

在以下研究中, 分别使用最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法和组合法对标普 500 指数收益率进行预测, 并对得到的 $R_{OS}^2$ 统计量的值展开分析。同时为了展示机器学习方法在预测方面的效果, 本文引入了最小二乘线性回归 (OLS) 作为另一个基准对比方法。

#### 4.3.1 未小波去噪

在未进行小波去噪的情况下, 各个预测算法的样本外预测效果如表 2 所示: 最小绝对收缩和选择算法的 $R_{OS}^2$ 值为-0.39%, CW-t 统计量结果为 1.17, 统计量不显著, 无法拒绝原假设。弹性网络算法的 $R_{OS}^2$ 值为 0.47%, CW-t 统计量结果为 1.95, 检验结果在显著性 0.05 的水平上显著不为零。主成分分析算法的 $R_{OS}^2$ 值为 0.19%, CW-t 统计量结果为 0.91, 统计量不显著, 无法拒绝原假设。缩放主成分分析算法的  $RR_{OS}^2$  值为 0.20%, CW-t 统计量结果为 2.13, 检验结果在显著性 0.05 的水平下显著不为零。偏最小二乘算法的 $R_{OS}^2$ 值为-0.04%, CW-t 统计量结果为 2.03, 检验结果在显著性 0.05 的水平下显著不为零。组合预测算法的 $R_{OS}^2$ 值为 1.30%, CW-t 统计量结果为 2.5, 检验结果在显著性 0.01 的水平下显著不为零。

表 2 小波去噪前后样本外预测效果<sup>1</sup>

预测方法	未去噪音		小波去噪	
	$R_{0s}^2$	CW-t	$R_{0s}^2$	CW-t
OLS	-9.05%	1.22	-3.05%	1.46*
lasso	-0.39%	1.17	1.76%	2.93***
enet	0.47%	1.95**	1.94%	2.97***
PCA	0.19%	0.91	0.69%	1.67**
SPCA	0.20%	2.13**	2.36%	3.3***
PLS	-0.04%	2.03**	2.27%	2.52***
COMB	1.30%	2.5***	2.77%	3.31***

比较各个预测算法的样本外预测结果可以发现，组合模型具有出色的样本外预测能力。在六种预测方法预测得到的 $R_{0s}^2$ 的值当中，最小绝对收缩和选择算法和偏最小二乘法进行预测得到的 $R_{0s}^2$ 值小于 0，而弹性网络、主成分分析法、缩放主成分分析、组合方法进行预测得到的 $R_{0s}^2$ 值大于 0。这表示最小绝对收缩和选择算法和偏最小二乘法的预测效果劣于现行平均基准方法，而弹性网络、主成分分析法、缩放主成分分析、组合方法的预测效果优于现行平均基准方法。其中，组合方法得到的 $R_{0s}^2$ 的值为 1.30%，明显大于其余五种单个预测方法所得到的 $R_{0s}^2$ 的值。考虑到月度市场收益率中难以预测的因素，月度 $R_{0s}^2$ 统计量超过 0.5%就足以产生显著的经济价值（Campbell 和 Thompson, 2008），因此这些结果充分说明了组合模型的预测能力。

比较各个预测算法的 CW-t 检验结果，可以发现组合方法得到的 $R_{0s}^2$ 的值体现出其预测能力相较于历史平均的基准方法有最为显著的提升。为检验 $R_{0s}^2$ 统计量是否显著大于 0，对 $R_{0s}^2$ 进行 CW-t 检验，原假设为 $R_{0s}^2 \leq 0$ ，相对于备择假设 $R_{0s}^2 > 0$ 。 $R_{0s}^2$ 显著大于 0

<sup>1</sup> \*、\*\*、\*\*\*分别代表 10%、5%、1%的显著性水平

即意味着该预测方法计算得到的均方预测误差显著小于历史平均的基准方法。根据假设检验结果，可以认为组合方法预测得到的 $R_{OS}^2$ 在显著性 0.05 的水平上显著；弹性网络、缩放主成分分析、偏最小二乘方法预测得到的 $R_{OS}^2$ 在 0.05 的水平上显著；而最小绝对收缩和选择算法和主成分分析算法得到的 $R_{OS}^2$ 不能认为显著大于 0，也即这两种算法的效果不能认为相较于历史平均的基准方法在对股票市场指数超过无风险利率的收益率的预测效果产生显著提升。

总结而言，组合预测方法的预测效果相较于历史平均的基准方法，在预测效果上的提升最为明显。组合预测算法的 $R_{OS}^2$ 值为 1.30%，在所有预测算法中 $R_{OS}^2$ 值的大小排名第一，且是排名第二的弹性网络算法的约三倍。CW-t 统计量结果为 2.5，检验结果在显著性 0.01 的水平下显著不为零，是所有预测算法中最为显著不为 0 的。这体现出了组合预测方法通过降低单一模型因为噪声和数据不足等原因而导致的错误率，减少误差和偏差以提高整体的鲁棒性的优势。

#### 4.3.2 小波去噪后

经过小波去噪处理后，各个预测算法的样本外预测效果得到了显著提升（见表 2）：最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 1.76%，CW-t 统计量结果为 2.93，检验结果在显著性 0.01 的水平下显著不为零。弹性网络算法的 $R_{OS}^2$ 值为 1.94%，CW-t 统计量结果为 2.97，检验结果在显著性 0.01 的水平下显著不为零。主成分分析算法的 $R_{OS}^2$ 值为 0.69%，CW-t 统计量结果为 1.67，检验结果在显著性 0.05 的水平上显著不为零。缩放主成分分析算法的 $R_{OS}^2$ 值为 2.36%，CW-t 统计量结果为 3.3，检验结果在显著性 0.01 的水平下显著不为零。偏最小二乘算法的 $R_{OS}^2$ 值为 2.27%，CW-t 统计量结果为 2.52，检验结果在显著性

0.01 的水平下显著不为零。组合预测算法的 $R_{OS}^2$ 值为 2.77%，CW-t 统计量结果为 3.31，检验结果在显著性 0.01 的水平下显著不为零。

组合方法仍然是六种方法中样本外预测效果最佳的。可以发现，六种预测模型预测得到的 $R_{OS}^2$ 的值当中，最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法和组合方法的 $R_{OS}^2$ 的值均大于零，且各算法相较于小波去噪处理以前的 $R_{OS}^2$ 的值均呈现出显著提升，考虑到月度市场收益率中难以预测的因素，月度 $R_{OS}^2$ 统计量超过 0.5%就足以产生显著的经济价值（Campbell 和 Thompson，2008），因此这些结果充分说明了各个预测模型基于小波去噪处理之后，均具有很好的预测效果。在六种预测模型中，组合预测算法的 $R_{OS}^2$ 的值达到 2.77，是所有预测算法当中最大的，且经过小波去噪处理后的组合预测算法的 $R_{OS}^2$ 的值是小波去噪处理前的组合预测算法的 $R_{OS}^2$ 的值的约 2 倍，再次验证了小波去噪处理对于提升预测模型预测效果的出色帮助。

经过小波去噪后，最小绝对收缩和选择算法、主成分分析算法弹性网络算法、缩放主成分分析与偏最小二乘法的检验结果显著性增强。对 $R_{OS}^2$ 进行 CW-t 检验，根据假设检验结果，发现最小绝对收缩和选择算法、弹性网络算法、主成分分析算法、缩放主成分分析、偏最小二乘法和组合方法所得到的 $R_{OS}^2$ 在不同显著性水平上显著。主成分分析法所得到的 $R_{OS}^2$ 在 0.05 的水平上显著；其余五种预测模型所得到的 $R_{OS}^2$ 均在 0.01 的水平上显著。比较小波去噪前后六种预测算法的 CW-t 检验结果，可以发现最小绝对收缩和选择算法与主成分分析算法由检验结果不显著转变为分别在 0.01 与 0.05 的显著性水平上显著，而弹性网络算法与缩放主成分分析、偏最小二乘法的检验结果也由在 0.05 的显著性水平上显著进一步转变为在 0.01 的显著性水平上显著。经过小波去噪处理后，六种预测算法在预

测效果上得到了提升，均在更大程度上显著优于历史均值的基准算法。

对 1978 至 2020 年间的递归估计窗口中出现频率最高的 30 个变量进行统计。其中，出现频率最高的 30 个自变量如表 3 所示。

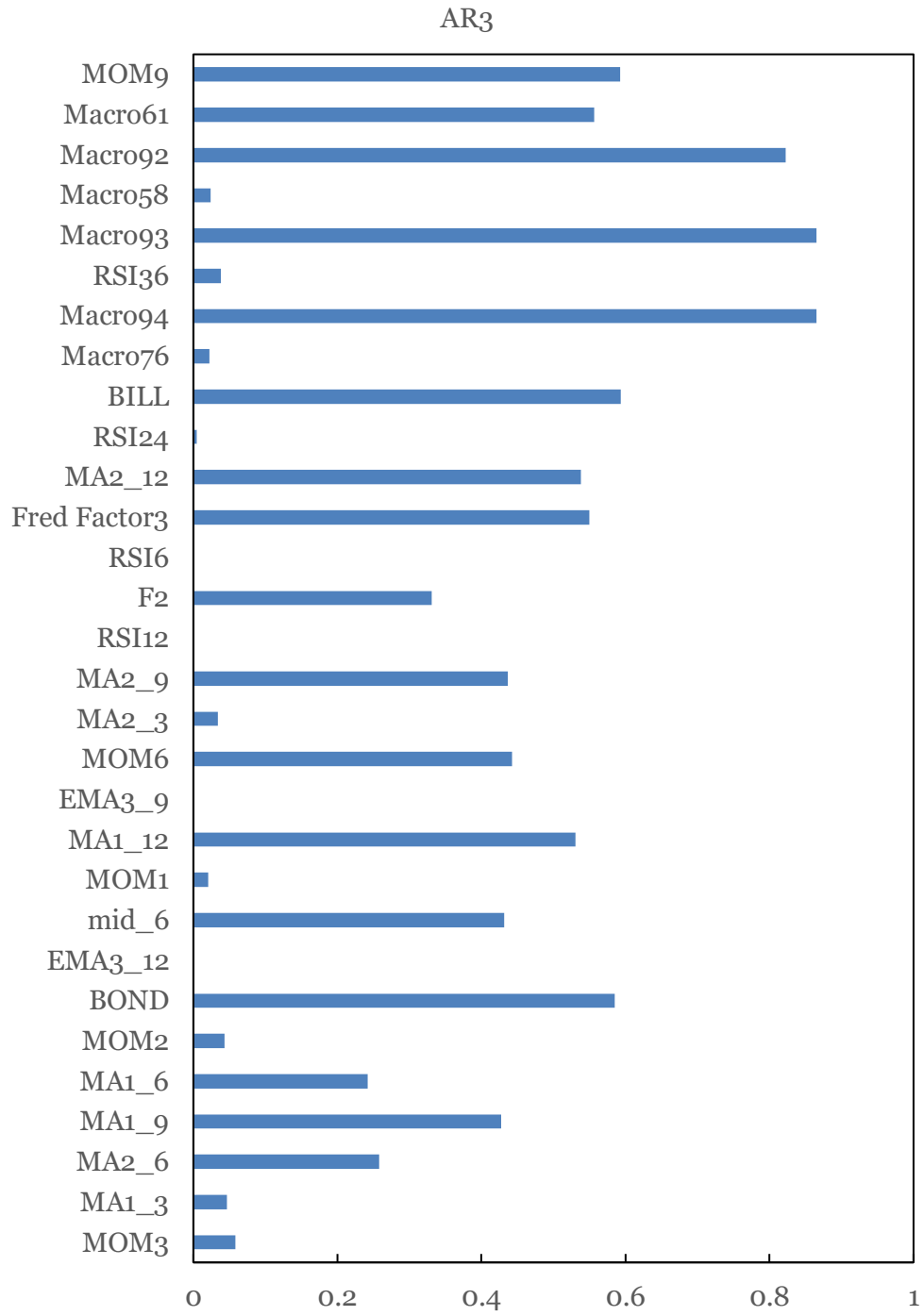
表 3 出现频率最高的 30 个变量介绍

序号	变量简称	变量解释
1	MOM9	9 个月动量指标
2	Macro61	耐用品的未成交订单
3	Macro92	10 年国债 C 减去联邦基金
4	Macro58	消费品新订单
5	Macro93	穆迪 Aaa 公司债券收益率减去联邦基金
6	RSI36	36 个月相对强弱指标
7	Macro94	穆迪 Baa 公司债券收益率减去联邦基金
8	Macro76	标准普尔综合普通股：股息收益
9	BILL	国库券利率减去其过去 12 个月的均值
10	RSI24	24 个月相对强弱指标
11	MA2_12	MA(2, 12)技术信号
12	Fred Factor3	对美国长期债券收益率分解得到的分量
13	RSI6	6 个月相对强弱指标
14	F2	对美国长期债券收益率分解得到的分量
15	RSI12	12 个月相对强弱指标
16	MA2_9	MA(2, 9)技术信号
17	MA2_3	MA(2, 3)技术信号
18	MOM6	6 个月动量指标
19	EMA3_9	EMA(3, 9)指数平均数指标
20	MA1_12	MA(1, 12)技术信号
21	MOM1	1 个月动量指标
22	mid_6	中期选举指标
23	EMA3_12	EMA(3, 12)指数平均数指标
24	BOND	长期政府债券利率减去其过去 12 个月的均值

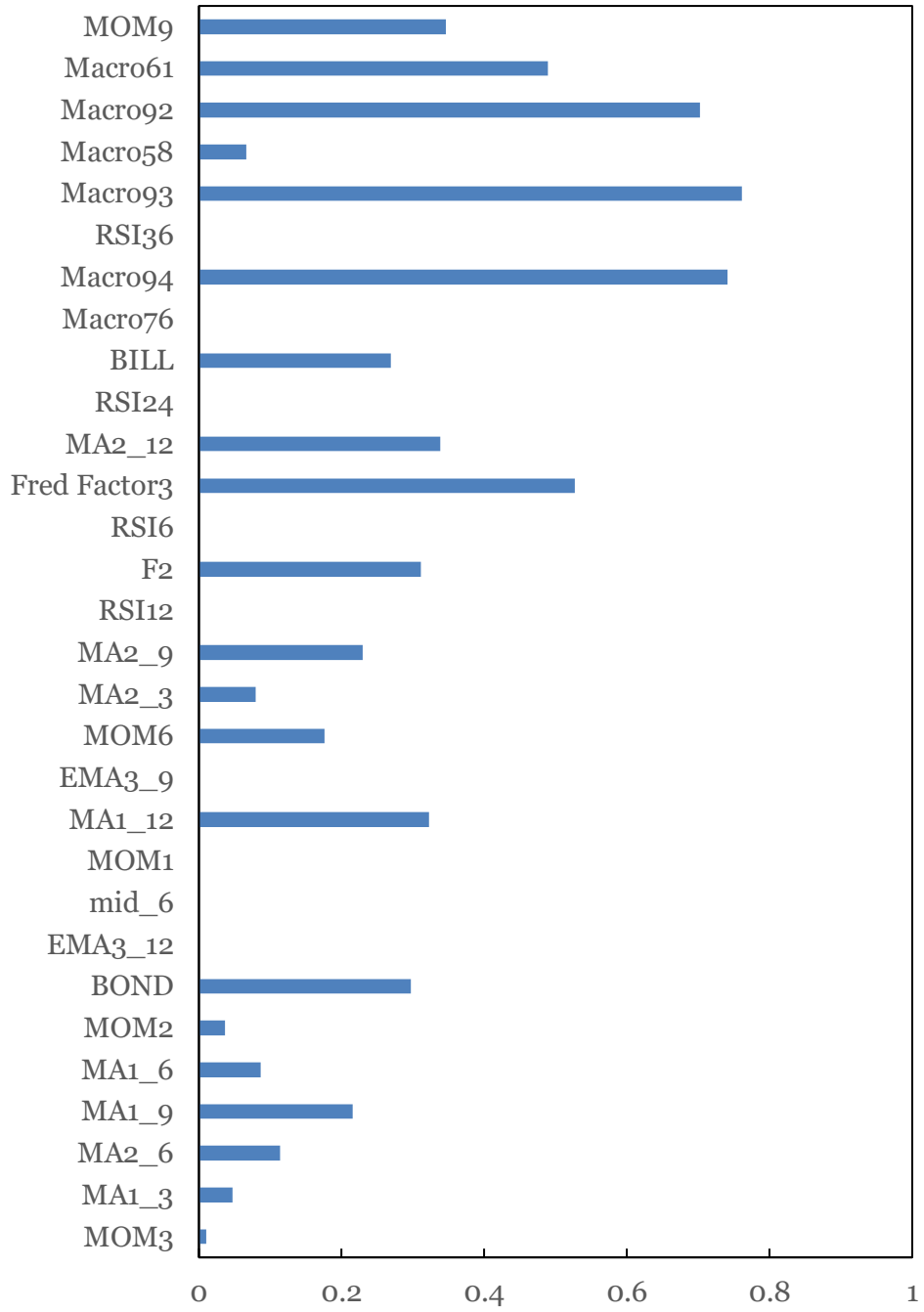
序号	变量简称	变量解释
25	MOM2	2 个月动量指标
26	MA1_6	MA(1, 6)技术信号
27	MA1_9	MA(1, 9)技术信号
28	MA2_6	MA(2, 6)技术信号
29	MA1_3	MA(1, 3)技术信号
30	MOM3	3 个月动量指标

进一步对这 30 个变量进行自回归分析。在本文中，观察到在 3 阶自回归模型中，所有变量都表现出一定程度的自相关性。这表明当前值与其过去的值存在一定的关联性。值得注意的是，宏观经济指标的自相关性普遍高于技术指标，这可能反映了宏观指标如 GDP、通胀率、利率等反映的经济条件具有更长期的连续性和稳定性。随着自回归阶数的增加至 6 阶，所有变量的自相关性开始减弱。这可能是因为随着时间跨度的增加，过去的的数据对当前值的影响减少。在这种情况下，市场的随机性和新出现的信息开始对价格产生更大的影响。当自回归阶数进一步增至 12 阶时，技术指标的自相关性几乎降至零。这反映了技术指标，如股价和交易量等，主要捕捉相对短期市场动态，而在长期数据中它们的相关性会显著减弱。相比之下，宏观经济指标即使在 12 阶自回归中仍保持一定程度的自相关性。这表明宏观经济指标反映的经济趋势和条件在更长的时间跨度内保持连续性。

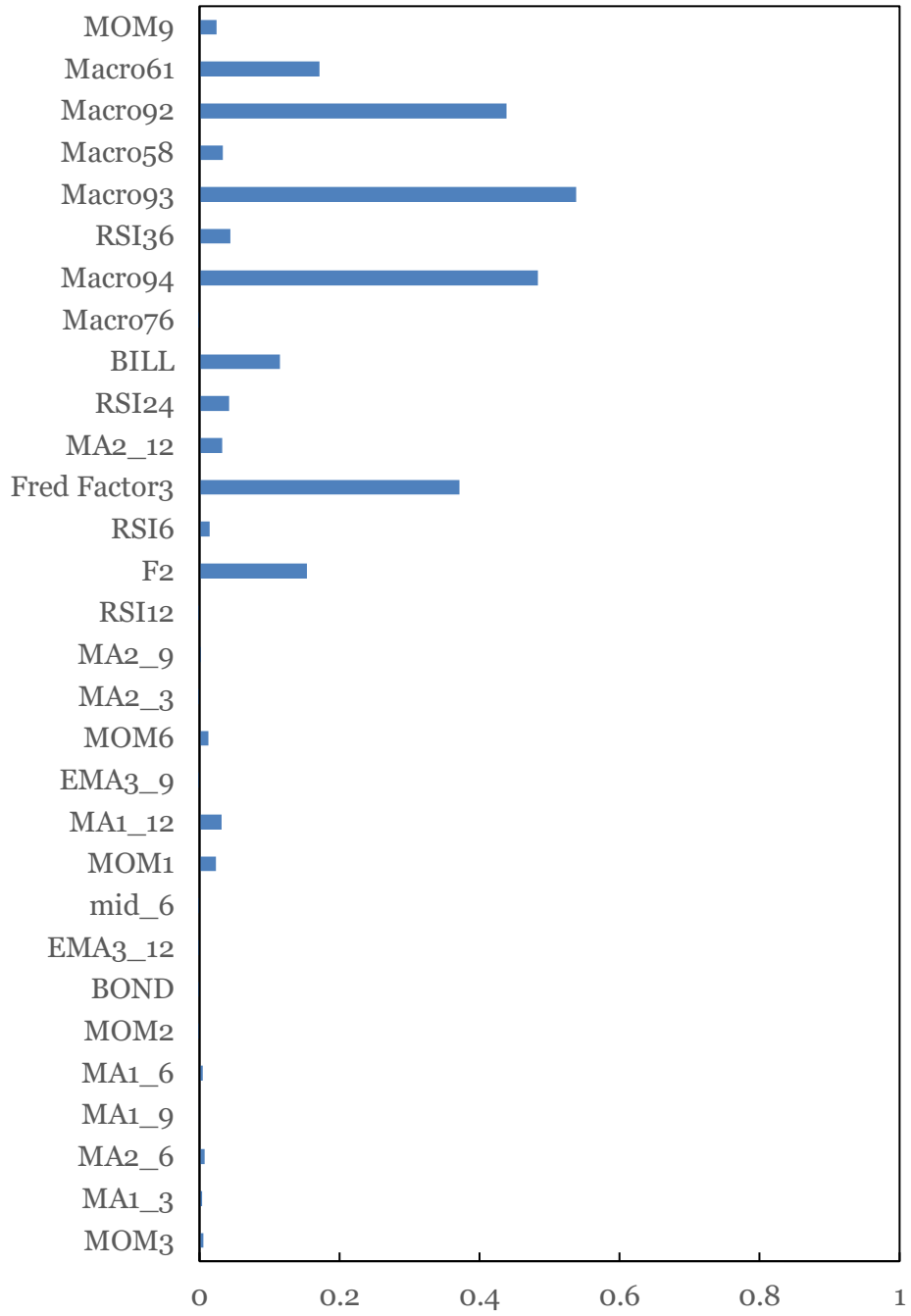
图 2 变量自回归结果



AR6



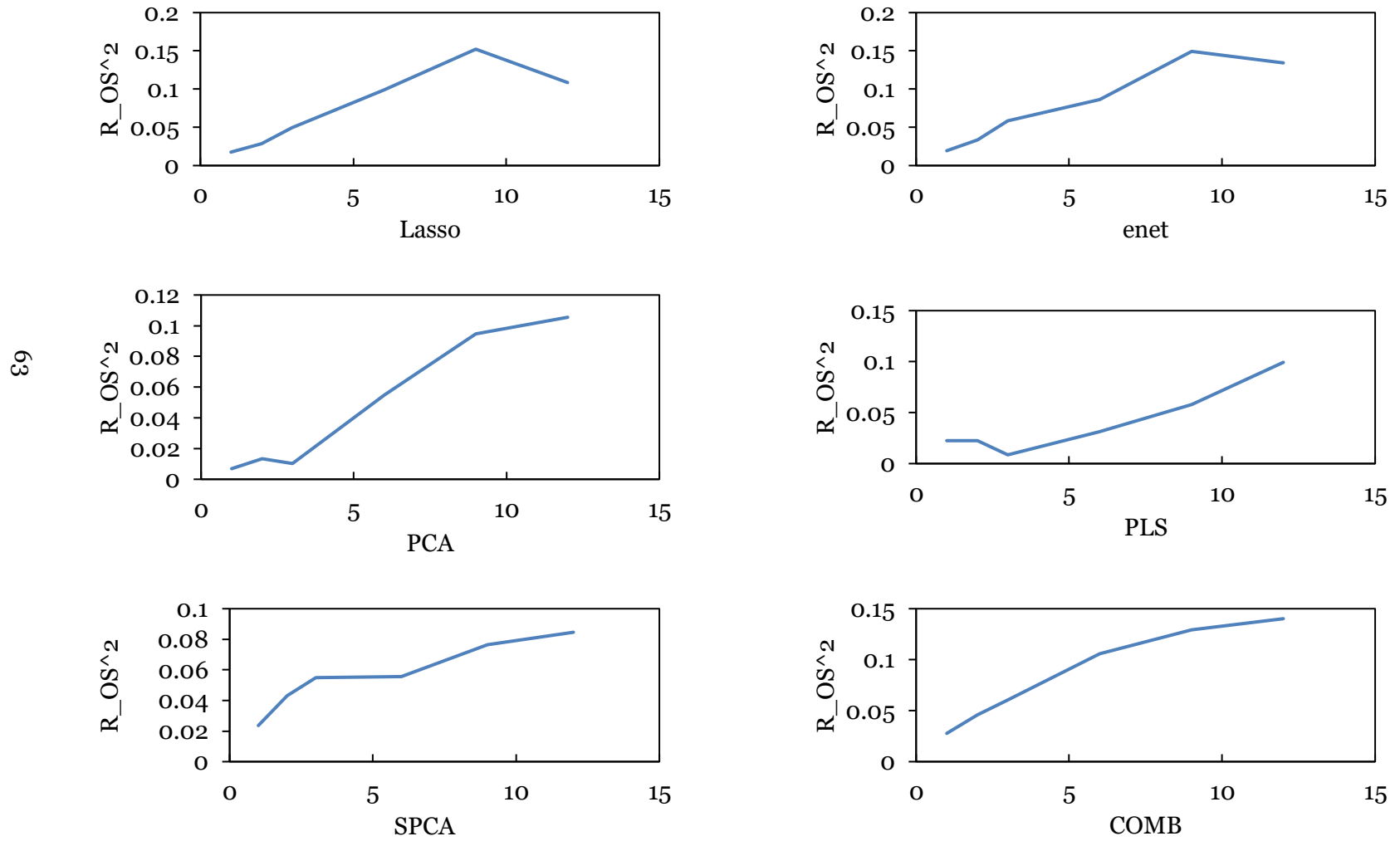
AR12



### 4.3.3 不同预测期限对预测效果的影响

使用月度宏观经济指标对未来 1-12 个月标普 500 指数收益进行预测，由图 3 可知，预测效果随着期限的改变而发生变化。从图中可以看出，最小绝对收缩和选择算法在某些时间点的预测性能有明显的下降；弹性网络算法的性能波动较大，但总体上呈上升趋势；主成分分析算法的曲线表明，随着预测时期的延长，模型的预测能力稳步提高；偏最小二乘算法的预测效果显示出一致的上升趋势；缩放主成分分析算法在中期预测时表现更好。组合算法的曲线相对平滑，并且随着预测时期的增加，其预测性能显著提高。

图 3 不同预测期限对预测效果的影响



整体来说，不同模型的性能差异可能反映了它们在处理特定数据集和预测时间框架时的优劣。一般来说，随着预测时间的延长，模型的不确定性可能会增加，但组合算法似乎在长期预测中表现较好，这表明模型组合可能提供了更稳定和可靠的预测结果。

## 五、投资组合实践

### 5.1 衡量指标

#### 5.1.1 年化收益率 (Annualized rate of return)

年化收益率是指将一项投资在一年内所获得的实际收益转换成百分比形式的指标。它是评估投资业绩的一个重要指标。年化收益率可以帮助投资者评估长期投资的表现。对于长期投资，实际收益率可能会因为各种因素而发生变化，但是年化收益率可以把这些波动统一到一年的时间框架内进行比较。

年化收益率的计算公式如下：

$$r = r_{monthly} \times 12$$

其中， $r$ 为年化收益率， $r_{monthly}$ 为月度收益率。

#### 5.1.2 确定性等价收益 (Certainty equivalent return)

确定性等价收益是一种用于衡量金融投资回报的指标，即风险调整后的收益率。它指的是投资者在相同的风险水平下可以获得的等价回报率，是为使无风险投资与风险投资有同样的吸引力而确定的无风险投资报酬率。也就是说，假设投资者有两种不同的投资选择，那么确定性等价收益率就是使得两种投资方案具有相同价值的回报率，也即在确定收益相同的情况下，能够提供与正在考虑的投资组合具有相同的效用值的收益率。

确定性等价收益率的计算通常需要考虑多种因素，例如投资的期限、风险水平、税收政策等。对于长期投资者而言，确定性等价收益率可以用来比较不同的投资方案，帮助他们做出更为理性的投资决策。计算公式如下：

$$CER_p = (\bar{R}_p - 0.5\gamma\sigma_p^2) \times 12$$

其中 $\bar{R}_p$ 和 $\sigma_p^2$ 是样本外实现的投资组合收益率的均值和方差， $\gamma$ 为风险厌恶系数。参照历史研究，本文将 $\gamma$ 值设定为3。

### 5.1.3 夏普比率（Sharpe ratio）

夏普比率是一种用来评估投资组合风险调整后收益率的指标，它通过将一個投资组合的超额收益率除以其标准差来计算。其中，超额收益率是指投资组合的年化收益率减去无风险收益率，标准差则代表着投资组合收益率的波动性。夏普比率是一种描述资产或投资组合的绩效指标，其反映了投资者每增加一个波动单位，收益率增加的值。夏普比率越高，表示投资组合在单位风险下创造的超额收益越多。一般来说，夏普比率越高，则意味着该投资组合的风险调整后收益率越高。月度夏普比率的计算公式如下：

$$\text{Monthly Sharpe ratio} = \frac{R_p - R_f}{\sigma_p}$$

其中 $R_p$ 为投资组合月度收益率， $R_f$ 为月度无风险利率， $\sigma_p$ 为投资组合月度收益率的标准差。为便于比较，将月度夏普比率转化为年度夏普比例，计算公式如下：

$$\text{Yearly Sharpe ratio} = \text{monthly Sharpe ratio} \times \sqrt{12}$$

### 5.1.4 最大回撤（Max drawdown）

最大回撤（Maximum Drawdown）是投资组合中的一个重要概念，可以用来评估投资组合的风险水平。最大回撤是指从一个投资组合净值曲线的高点到低点的最大跌幅，也即投资组合遭受的最大损失。最大回撤是一个重要的风险指标，因为它可以告诉投资者，如果他们在投资组合最高点时购买该投资组合，他们可能会面临的最大的亏损，并且帮助他们决定是否应该持有或退出某个投资组合。最大回撤的数值越大，表示投资组合的风险越高。假设一个投资组合的净值曲线是 $P$ ，最大回撤的计算公式如下：

$$MDD = \max \left[ \frac{(P_i - P_j)}{P_i} \right]$$

其中， $P_i$ 是任意时刻的净值， $P_j$ 是该时刻之前的最高净值。最大回撤的计算方法表示从任意时刻到前面的最高点之间的最大跌幅。

## 5.2 手续费计算

本文中，假设股票的手续费率为 50 个基点，也就是说每笔交易的手续费为交易金额的 5%。这是一种常见的股票交易费用设置方式，通常由交易所或经纪公司制定，并根据交易规模和交易类型进行调整。手续费是指交易时需要支付给经纪公司或交易所的费用，用于维护交易平台和提供服务。在交易股票时，手续费是必须要支付的费用之一，而且通常是与交易金额成比例的。因此，在进行股票交易时，需要考虑手续费对交易成本的影响，以便做出明智的决策。

实际上，标普 500 交易时的手续费远低于 50 个基点，本文选择了一个较高的手续费用，以此来模拟交易时的冲击成本。因为标普 500 ETF 的流动性非常好，即使是非常大的资金，其冲击成本通常也很小，所以 50 个基点的双向手续费是对从理论到实盘交易的一种非常好的模拟。

## 5.3 择时策略

### 5.3.1 策略介绍

择时策略是指根据市场走势和投资者的判断，在合适的时机进行买入或卖出的决策。风险调整后收益是指在考虑风险因素的情况下，通过投资获得的收益。风险因素包括市场风险、行业风险、公司风险等。在进行标普 500 指数的交易时，需要考虑到市场的整体风险，以及相关公司的风险和机会。

择时策略的具体实现方式，是通过求解目标函数以最大化投资者效用：

$$\arg \max_{w_{t+1|t}} (w_{t+1|t} \hat{r}_{t+1|t} - 0.5\gamma w_{t+1|t}^2 \hat{\sigma}_{t+1|t}^2)$$

该目标函数的解为：

$$w_{t+1|t}^* = \left(\frac{1}{\gamma}\right) \left(\frac{\hat{r}_{t+1|t}}{\hat{\sigma}_{t+1|t}^2}\right)$$

其中， $\hat{r}_{t+1|t}$ 为预测得到的超额收益； $\hat{\sigma}_{t+1|t}^2$ 为基于过去 60 个月的超额收益率序列计算得到的方差，也即风险； $\gamma$ 为相对风险厌恶系数。假定投资者有一定的风险厌恶程度，设置 $\gamma = 3$ ，并将投资者在 $t + 1$ 时刻的持仓值 $w_{t+1|t}$ 限制在 $[0, 150\%]$ 区间内。根据求解得到的 $w_{t+1|t}^*$ 动态调整标普 500 指数的持仓情况，对标普 500 指数进行投资。

### 5.3.2 策略回测结果

#### 1) 不考虑交易手续费

表 4 展示了使用最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法、组合方法、历史均值模型进行预测回归的投资组合的表现。其中历史均值模型为基准模型。同时为了展示机器学习方法在预测方面的效果，本文引入了最小二乘线性回归（OLS）作为另一个基准对比方法。本文根据以上择时策略的构建方法，根据对市场超额收益的预测，选择增加或减少对标普 500 指数基金的持仓，从而构建了一个择时投资组合。在没有使用小波分析对股票收益时间序列进行去噪时，最小绝对收缩和选择算法的年化收益率为 11.11%，年化确定性等价收益率为 7.80%，年化夏普比率为 0.47；弹性网络算法的年化收益率为 12.62%，年化确定性等价收益率为 9.24%，年化夏普比率为 0.56；主成分分析法的年化收益率为 10.34%，年化确定性等价收益率为 7.86%，年化夏普比率为 0.48；缩放主成分分析的年化收益率为 12.62%，年化确定性等价收益率为

9.64%，年化夏普比率为 0.6；偏最小二乘法的年化收益率为 10.42%，年化确定性等价收益率为 6.91%，年化夏普比率为 0.41；组合方法的年化收益率为 12.15%，年化确定性等价收益率为 9.09%，年化夏普比率为 0.56；历史均值法年化收益率为 8.31%，年化确定性等价收益率为 6.22%，年化夏普比率为 0.35。

表 4 基于各预测模型的投资组合策略收益

预测方法	未去噪音			小波去噪		
	return	cer	sharpe ratio	return	cer	sharpe ratio
OLS	10.77%	8.00%	0.49	10.68%	8.34%	0.52
lasso	11.11%	7.80%	0.47	13.97%	11.22%	0.73
enet	12.62%	9.24%	0.56	13.87%	11.15%	0.72
PCA	10.34%	7.86%	0.48	12.13%	8.94%	0.54
SPCA	12.62%	9.64%	0.6	14.81%	11.74%	0.74
PLS	10.42%	6.91%	0.41	12.28%	9.41%	0.59
COMB	12.15%	9.09%	0.56	13.91%	11.18%	0.72
hmean	8.31%	6.22%	0.35	8.31%	6.22%	0.35

可以看到，最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法和组合方法六种预测模型下构建的投资组合对应的年化收益率和年化等价性确定收益率均高于历史均值模型的对应指标。其中，缩放主成分分析算法的年化收益率和年化等价性确定收益率的值在构建的所有投资组合中值最大，体现了其较好的收益回报，其次是弹性网络算法与组合分析方法。可以发现，最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法和组合方法六种预测模型下构建的投资组合对应的夏普比率均高于历史均值模型的对应指标，其中对应夏普比率最高的是缩放主成分分析算法，值为 0.6，代表着投资者每多增加 1 单位波动，能够额外得到 0.6 单位的收益。以夏普比率作为排序依据，排在第二位的是弹性网络算法与组合方法，值为 0.56。注

意到组合方法的年化收益率和年化等价性确定收益率分别为 12.15%和 9.09%，略低于弹性网络算法，而弹性网络算法和组合方法的夏普比率是相等的。因此可以确认组合方法在减小投资组合波动率方面相较其他算法有着一定的优势。

使用小波分析技术对股票收益数据进行去噪音后，本文发现：最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法和组合方法六种预测模型下构建的投资组合对应的年化收益率和年化等价性确定收益率均高于历史均值模型的对应指标。这六种预测模型基于小波处理后的数据构建的投资组合，与同一预测模型基于原始数据构建的投资组合相比，在年化收益率和年化等价性确定收益率上均有一定程度的提升。由此可以发现小波去噪处理对于构建更优的投资组合有显著效果。其中，缩放主成分分析算法的年化收益率和年化等价性确定收益率的值在构建的所有投资组合中值最大，年化收益率达到了 14.81%，年化等价性确定收益率达到了 11.74%，体现了其较好的收益回报。最小绝对收缩和选择算法、弹性网络算法与组合分析方法的年化收益率均达到了 13%以上，年化等价性确定收益率均达到了 11%以上。相较于历史均值模型 8.31%的年化收益率和 6.22%的年化等价性确定收益率，基于缩放主成分分析算法、最小绝对收缩和选择算法、弹性网络算法与组合分析方法生成的投资组合在投资回报上有显著的提升。

同时，缩放主成分分析算法、最小绝对收缩和选择算法、弹性网络算法与组合分析方法四种预测模型下构建的投资组合对应的夏普比率均高于 0.7。这一指标的差异体现了缩放主成分分析算法、最小绝对收缩和选择算法、弹性网络算法与组合分析方法四种算法相较于历史均值模型在降低投资组合波动率有着显著优势。在承受相同单位风险的情况下，使用以上四种预测算法构建投资组合能够给投资者带来更高的收益。

以最小二乘线性回归作为比较基准，可以发现包括最小绝对收缩和选择算法、主成分分析法、和弹性网络在内的机器学习方法在对标普 500 指数的预测上相较于最小二乘回归都有显著提升。这是因为这些机器学习方法具备处理更复杂数据关系的能力。不同于最小二乘线性回归仅能处理线性关系，这些机器学习模型能够识别和利用数据中的非线性模式和隐藏的关系，有效地处理高维数据集中的复杂关系。这在金融市场预测中尤为重要，因为市场数据通常包含多种因素和内在的复杂性。因此，相比于传统的最小二乘线性回归方法，这些机器学习技术在预测标普 500 指数方面表现出更高的准确性和效率。

## 2) 考虑交易手续费

表 5 展示了考虑交易手续费和冲击成本后，分别使用最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法、组合方法、最小二乘法、历史均值模型进行投资所获得的收益率、确定性等价收益和夏普比率。其中历史均值模型为基准模型。

表 5 考虑交易手续费和冲击成本后的投资组合收益情况

预测方法	未去噪音			小波去噪		
	50bps_ return	50bps_ cer	50bps_ sharpe ratio	50bps_ return	50bps_ cer	50bps_ sharpe ratio
lasso	9.02%	5.71%	0.33	12.56%	9.80%	0.62
enet	10.35%	6.97%	0.41	12.43%	9.69%	0.61
PCA	9.76%	7.29%	0.44	11.55%	8.36%	0.51
SPCA	11.42%	8.44%	0.51	14.00%	10.92%	0.69
PLS	9.38%	5.86%	0.34	9.12%	6.24%	0.36
COMB	10.63%	7.58%	0.45	12.59%	9.85%	0.62
OLS	7.00%	4.23%	0.21	7.95%	5.58%	0.3
hmean	8.14%	6.06%	0.34	8.14%	6.06%	0.34

当将 50 个基点的交易手续费纳入考虑范围后，最小绝对收缩和选择算法的年化收益

率为 9.02%，年化确定性等价收益率为 5.71%，年化夏普比率为 0.33；弹性网络算法的年化收益率为 10.35%，年化确定性等价收益率为 6.97%，年化夏普比率为 0.41；主成分分析法的年化收益率为 9.76%，年化确定性等价收益率为 7.29%，年化夏普比率为 0.44；缩放主成分分析的年化收益率为 11.42%，年化确定性等价收益率为 8.44%，年化夏普比率为 0.51；偏最小二乘法的年化收益率为 9.38%，年化确定性等价收益率为 5.89%，年化夏普比率为 0.34；组合方法的年化收益率为 10.63%，年化确定性等价收益率为 7.58%，年化夏普比率为 0.45；历史均值模型年化收益率为 8.14%，年化确定性等价收益率为 6.06%，年化夏普比率为 0.34。

由于考虑了交易手续费，增加了交易成本，基于最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法、组合方法六种预测模型构建的投资组合对应的年化收益率、年化确定性等价收益率和夏普比率均有一定程度的下降。其中最小绝对收缩和选择算法和弹性网络算法受交易成本的影响最大，其对应的年化收益率与年化确定性等价收益率分别下降了约 2%，说明这两种算法对应的投资组合发生股票仓位调整的频率较高。相较而言，缩放主成分分析法、偏最小二乘法、组合方法的年化收益率和年化确定性等价收益率受交易成本的影响较小，均在 1%左右。主成分分析算法的年化收益率和年化确定性等价收益率变化仅在 0.5%左右。这六种预测算法的夏普比率也分别有一定程度的下降，但仍明显高于历史均值模型的夏普比率。其中，缩放主成分分析算法的夏普比率最高，为 0.51；组合方法的夏普比率第二高，为 0.45。这意味着投资者如果使用缩放主成分分析算法或是组合方法，每多承担一单位风险，相较使用夏普比率为 0.34 的历史均值模型，可以多获得约 0.1~0.15 单位的收益。

当将 50 个基点的交易手续费纳入考虑范围后，使用小波去噪技术构建的投资组合的结果如下：得到最小绝对收缩和选择算法的年化收益率为 12.56%，年化确定性等价收益率为 9.80%，年化夏普比率为 0.62；弹性网络算法的年化收益率为 12.43%，年化确定性等价收益率为 9.69%，年化夏普比率为 0.61；主成分分析法的年化收益率为 11.55%，年化确定性等价收益率为 8.36%，年化夏普比率为 0.51；缩放主成分分析的年化收益率为 14.00%，年化确定性等价收益率为 10.92%，年化夏普比率为 0.69；偏最小二乘法的年化收益率为 9.12%，年化确定性等价收益率为 6.24%，年化夏普比率为 0.36；组合方法的年化收益率为 12.59%，年化确定性等价收益率为 9.85%，年化夏普比率为 0.62。

由于考虑了交易手续费，增加了交易成本，基于最小绝对收缩和选择算法、弹性网络、主成分分析法、缩放主成分分析、偏最小二乘法、组合方法六种预测模型构建的投资组合对应的年化收益率、年化确定性等价收益率和夏普比率均有一定程度的下降。其中偏最小二乘算法受交易成本的影响最大，其对应的年化收益率与年化确定性等价收益率分别下降了约 3%，说明该算法对应的投资组合发生股票仓位调整的频率较高。相较而言，最小绝对收缩和选择算法，弹性网络算法和组合方法的年化收益率和年化确定性等价收益率受交易成本的影响较小，均在 1%左右。主成分分析算法和缩放主成分分析算法的年化收益率和年化确定性等价收益率变化仅在 0.5%-0.8%。这六种预测算法的夏普比率也分别有一定程度的下降，但仍明显高于历史均值模型的夏普比率。其中，缩放主成分分析算法的夏普比率最高，为 0.69；最小绝对收缩和选择算法和组合方法的夏普比率第二高，为 0.62。这意味着投资者如果使用缩放主成分分析算法或是组合方法，每多承担一单位风险，可以额外获得历史均值模型对应收益的两倍。除此之外，本文也发现，这六种预测模

型基于小波处理后的数据构建的投资组合，与同一预测模型基于原始数据构建的投资组合相比，在年化收益率和年化等价性确定收益率上均有一定程度的提升。由此可以发现小波去噪处理对于构建更优的投资组合有显著效果。

### 5.3.3 变量类型对回测结果的影响

将 221 个变量按照变量类型划分为 45 个技术指标与 176 个宏观指标。使用小波去噪后的收益率数据，对技术指标组/宏观指标组分别进行变量筛选、样本外预测并应用择时策略，得到的回测结果如表 6 所示。

表 6 使用不同类型指标的策略回测结果<sup>1</sup>

技术加宏观								
	$R_{OS}^2$	CW-t	return	cer	sharpe ratio	50bps_return	50bps_cer	50bps_sharpe ratio
lasso	1.76%	2.93***	13.97%	11.22%	0.73	12.56%	9.80%	0.62
enet	1.94%	2.97***	13.87%	11.15%	0.72	12.43%	9.69%	0.61
PCA	0.69%	1.67**	12.13%	8.94%	0.54	11.55%	8.36%	0.51
PLS	2.27%	2.52***	12.28%	9.41%	0.59	9.12%	6.24%	0.36
SPCA	2.36%	3.3***	14.81%	11.74%	0.74	14.00%	10.92%	0.69
COMB	2.77%	3.31***	13.91%	11.18%	0.72	12.59%	9.85%	0.62
h_mean			8.31%	6.22%	0.35	8.14%	6.06%	0.34
技术								
	$R_{OS}^2$	CW-t	return	cer	sharpe ratio	50bps_return	50bps_cer	50bps_sharpe ratio
lasso	-0.38%	0.97	12.53%	9.65%	0.6	11.53%	8.66%	0.53
enet	-0.51%	0.85	12.27%	9.43%	0.59	11.22%	8.38%	0.51
PCA	0.27%	1.01	10.67%	8.10%	0.49	10.30%	7.73%	0.47
PLS	0.35%	1.95**	8.99%	7.20%	0.44	7.72%	5.93%	0.32
SPCA	0.33%	1.05	10.24%	7.82%	0.48	9.79%	7.37%	0.44

<sup>1</sup> \*、\*\*、\*\*\*分别代表 10%、5%、1%的显著性水平

COMB	0.33%	1.13	10.80%	8.51%	0.54	10.01%	7.72%	0.47
宏观								
	$R_{Os}^2$	CW-t	return	cer	sharpe ratio	50bps_r eturn	50bps_ cer	50bps_ sha rpe ratio
lasso	1.00%	2.30**	12.88%	9.80%	0.61	11.26%	8.17%	0.49
enet	1.40%	2.49***	13.15%	10.09%	0.63	11.54%	8.46%	0.51
PCA	1.56%	2.81***	12.83%	9.69%	0.6	11.56%	8.43%	0.51
PLS	2.18%	2.51***	12.15%	9.30%	0.58	9.03%	6.17%	0.35
SPCA	1.78%	2.79***	14.45%	10.46%	0.63	13.21%	9.20%	0.55
COMB	2.20%	2.87***	13.63%	10.51%	0.66	12.04%	8.91%	0.55

由表 6 可知，仅使用技术指标构建的策略在扣除 50 个基点的手续费之后，年化收益率约为 7%-12%，确定性等价收益约为 5%-9%，夏普比率约为 0.3-0.6；仅使用宏观指标构建的策略的年化收益率约为 9%-14%，确定性等价收益约为 6%-10%，夏普比率约为 0.3-0.6。相较而言，使用宏观指标构建的策略效果要优于技术指标构建的策略效果。这可能是由于宏观指标能够提供更加深入和全面的市场分析。消费指数、通货膨胀率、利率和政策变化等宏观指标，反映了整体经济和市场的健康状况，使策略能够更好地适应和预测市场的月频波动。相比之下，技术指标虽然能够有效捕捉市场的短期动态和趋势，但它们主要基于历史价格和交易量数据，可能不足以全面反映即将发生的宏观经济变化。因此，在这种需要预测中长期市场走势的场景中，宏观指标提供的信息可能更为关键和有效。

对比仅使用技术/宏观指标的策略回测结果与使用技术+宏观指标构建得到的策略回测结果，可以发现使用技术+宏观指标构建得到的策略在回测结果上表现最好。在扣除 50 个基点的手续费之后，年化收益率约为 9%-14%，确定性等价收益约为 6%-11%，夏普比率约为 0.3-0.7，略高于另外两种策略。这一现象表明，将技术指标与宏观指标相结合使用，在构建投资策略时可以提供更加全面和均衡的视角。技术指标和宏观指标各有其独特

的优势和局限性，通过将它们结合起来，可以互补彼此的不足，从而更有效地捕捉市场机会并降低风险。

#### 5.3.4 累积财富积累曲线

图 4 择时策略的收益率曲线

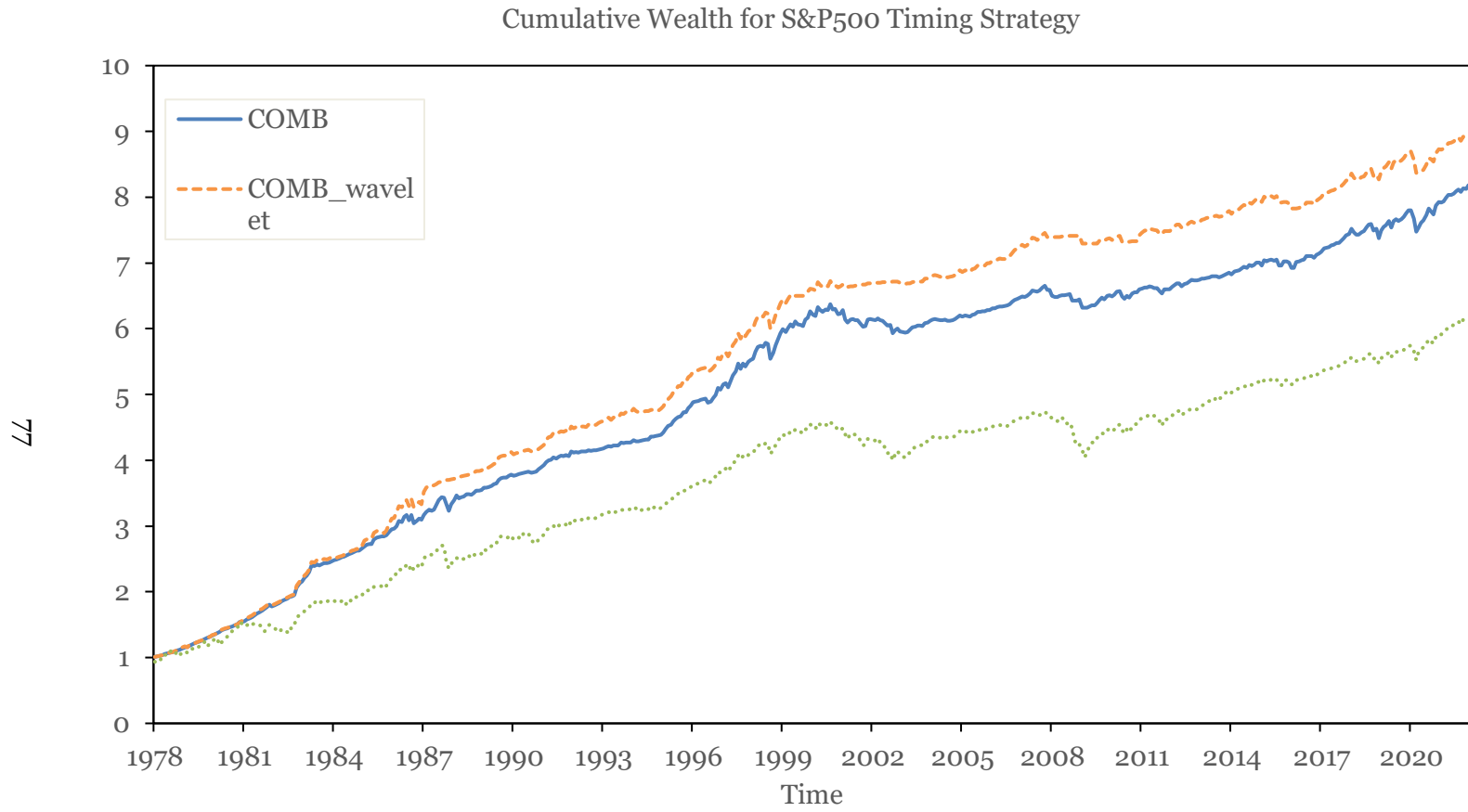


图 4 展示了按照 S&P500 指数收益率曲线（绿色）、按照组合预测方法进行投资获得的收益率曲线（蓝色）以及小波变换多分辨率分析处理后的组合预测方法进行投资所获得的收益率曲线（橙色）。

可以发现，使用组合预测方法的两条收益率曲线（蓝色、橙色）均位于 S&P500 指数收益率曲线（绿色）的上方。组合预测方法的曲线斜率更高，且相对 S&P500 指数要更为平稳。这意味着根据这组合预测方法投资可以有效预测经济扩张与衰退期间的股票月度收益率，比其他预测模型包含更多实际收益率的信息，这使得投资者可以避免一些可能在经济衰退时期的大幅度损失，从而能够收获更高的投资收益率。此外，经过去噪处理的组合预测方法所得到的收益率曲线（橙色）在仅使用组合预测方法的曲线（蓝色）上方。这意味着通过小波去噪处理可以帮助投资者提高预测精度，从而提高投资收益率。

## 六、选股策略

### 6.1 行业组合

在构建标普 500 收益率预测模型时，将股票按照行业划分是一项至关重要的任务。由于行业受到不同的宏观经济和行业特定因素的影响，股票表现具有高度异质性。因此，按照所属行业对股票进行分类，有助于更好地理解 and 解释收益率的波动和预测。此外，同一行业的股票通常会受到相似的宏观经济和行业特定因素的影响。例如，技术行业的股票通常会受到技术进步和市场需求等因素的影响，而金融行业的股票通常会受到利率和货币政策等因素的影响。因此，按照行业分类，不仅可以更好地理解 and 预测整个指数的表现，还可以更好地区分不同行业之间的差异，从而更好地预测股票的表现。

在本文中，将标普 500 指数中包含的所有股票根据其所属行业划分为“消费”、“制造”、“高科技”、“健康”和“其他”五个行业（表 7）。本文将基于此行业分类来评估各预测模型的预测效果以及构建的投资组合效果。

表 7 行业划分介绍

行业名称	股票分布数量 (1986.1)	股票分布数量 (2020.12)	市值占比 (1986.1)	市值占比 (2020.12)
消费	583	474	20.7%	17.9%
制造	844	548	48.8%	13%
高科技	184	623	18.2%	38%
健康	38	708	4.1%	9.3%
其他	339	1007	8.2%	21.8%

## 6.2 策略介绍

### 1) 行业轮动方法 (industry rotation)

核心思想是根据历史数据计算出各行业的预期收益率，并且买入预期收益率高于平均预期收益率行业的股票，卖出预期收益率低于平均预期收益率行业的股票，从而构建买入-卖空投资组合，这个组合的总持仓为 0，波动率等于市场组合的波动率，从而使得其收益率与其他投资组合是可比的。投资组合的具体构建方法如下：

$$w_{i,t} = C_t(\hat{r}_{i,t} - \bar{r}_t)$$

其中， $w_{i,t}$  为 t 时刻第 i 个行业的权重， $\hat{r}_{i,t}$  为 t 时刻第 i 个行业收益率的预测值， $\bar{r}_t$  为 t 时刻 5 个行业收益率的预测值的均值。 $C_t$  满足以下等式：

$$C_t = \frac{\hat{\sigma}_{mt}}{\hat{\sigma}_{pt}}$$

其中  $\hat{\sigma}_{mt}$  为标普 500 指数历史收益率序列的标准差， $\hat{\sigma}_{pt}$  为基于行业轮动方法构建的投资组合历史收益率序列的标准差。这种构建方式保证了行业轮动投资组合的市场风险敞口接近于 0，波动率接近于标普 500 指数的波动率。

行业轮动策略旨在通过合理分散投资风险来获取更高的收益。使用行业轮动策略的优点在于它可以使投资者根据市场情况灵活地调整投资组合，从而获得更好的回报。此外，该策略还可以有效地降低投资组合的风险。

### 2) 行业等权重方法 (industry equal weight)

将标普 500 指数划分为消费、制造、高科技、健康和其他共五个行业，对每个行业分配相同投资组合权重，即 20%。

### 3) 行业均值方差方法 (industry MV)

在 Markowitz (1952) 的均值-方差模型中, 投资者在投资组合回报的均值和方差之间进行权衡。要实现这个模型, 本文遵循经典的“插件”方法: 也就是说, 将等式:

$$\max_{\mathbf{w}_t} \left( \mathbf{w}_t^\top \boldsymbol{\mu}_t - \frac{\gamma}{2} \mathbf{w}_t^\top \boldsymbol{\Sigma}_t \mathbf{w}_t \right), \text{ s.t. } \mathbf{1}_N^\top \mathbf{w}_t = 1$$

当中资产收益均值 $\boldsymbol{\mu}_t$ 和协方差矩阵 $\boldsymbol{\Sigma}_t$ 分别替换为它们的样本收益均值 $\hat{\boldsymbol{\mu}}_t$ 和样本协方差矩阵 $\hat{\boldsymbol{\Sigma}}_t$ , 并求解 $\mathbf{w}_t$ 。其中 $\gamma$ 为投资者的风险厌恶系数,  $\mathbf{w}_t$ 为投资组合五个行业组的权重向量。需要注意的是, 这种投资组合策略忽略了可能存在的估计误差。

### 4) 行业最小方差方法 (industry GMV)

在行业最小方差方法下, 本文选择风险资产组合, 以最小化方差。具体目标函数与约束条件如下:

$$\min_{\mathbf{w}_t} \mathbf{w}_t^\top \boldsymbol{\Sigma}_t \mathbf{w}_t \text{ s.t. } \mathbf{1}_N^\top \mathbf{w}_t = 1$$

为了实现此策略, 本文仅使用协方差矩阵的估计值 (样本协方差矩阵 $\hat{\boldsymbol{\Sigma}}_t$ ) 并完全忽略对预期回报的估计值。 $\mathbf{w}_t$ 为投资组合五个行业组的权重向量。

## 6.3 行业样本外预测结果

在消费、制造、高科技、健康和其他共五个行业中分别应用六种预测模型进行基于小波去噪处理后的样本外预测, 得到以下结果 (表 8):

表 8 各行业样本外预测结果<sup>1</sup>

	Cnsmr		Manuf		HiTec		Hlth		Other	
	$R^2_{OS}$	CW-t	$R^2_{OS}$	CW-t	$R^2_{OS}$	CW-t	$R^2_{OS}$	CW-t	$R^2_{OS}$	CW-t
lasso	1.88%	3.54***	1.01%	2.44***	3.20%	3.98***	1.21%	2.38***	1.35%	2.31**
enet	2.44%	3.78***	1.13%	2.47***	3.33%	4.06***	1.48%	2.6***	2.05%	2.63***
PCA	2.10%	3.11***	1.80%	3.27***	1.41%	2.84***	0.50%	1.82**	0.13%	1.1
PLS	0.88%	2.03**	2.12%	2.4***	0.73%	1.5*	1.42%	2.86***	0.62%	1.9**
SPCA	3.38%	3.87***	1.71%	2.71***	1.68%	3.16***	1.44%	2.49***	1.44%	2.58***
COMB	3.83%	4.13***	2.28%	3.1***	2.94%	4.2***	1.88%	3.07***	2.19%	2.83***

<sup>1</sup> \*、\*\*、\*\*\*分别代表 10%、5%、1%的显著性水平

在消费行业，最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 1.88%，CW-t 检验统计量结果为 3.54，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；弹性网络算法的 $R_{OS}^2$ 值为 2.44%，CW-t 检验统计量结果为 3.78，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；主成分分析算法的 $R_{OS}^2$ 值为 2.10%，CW-t 检验统计量结果为 3.11，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；偏最小二乘算法的 $R_{OS}^2$ 值为 0.88%，CW-t 检验统计量结果为 2.03，认为在 0.05 的显著性水平上 $R_{OS}^2$ 显著大于 0；缩放主成分分析算法的 $R_{OS}^2$ 值为 3.38%，CW-t 检验统计量结果为 3.87，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；组合预测模型的 $R_{OS}^2$ 值为 3.83%，CW-t 检验统计量结果为 4.13，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0。

在制造行业，最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 1.01%，CW-t 检验统计量结果为 2.44，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；弹性网络算法的 $R_{OS}^2$ 值为 1.13%，CW-t 检验统计量结果为 2.47，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；主成分分析算法的 $R_{OS}^2$ 值为 1.80%，CW-t 检验统计量结果为 3.27，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；偏最小二乘算法的 $R_{OS}^2$ 值为 2.12%，CW-t 检验统计量结果为 2.4，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；缩放主成分分析算法的 $R_{OS}^2$ 值为 1.71%，CW-t 检验统计量结果为 2.71，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；组合预测模型的 $R_{OS}^2$ 值为 2.28%，CW-t 检验统计量结果为 3.1，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0。

在高科技行业，最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 3.20%，CW-t 检验统计量结果为 3.98，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；弹性网络算法的 $R_{OS}^2$ 值为 3.33%，CW-t 检验统计量结果为 4.06，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；主成分分析算法的 $R_{OS}^2$ 值为 1.41%，CW-t 检验统计量结果为 2.84，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显

著大于 0；偏最小二乘算法的 $R_{OS}^2$ 值为 0.73%，CW-t 检验统计量结果为 1.5，认为在 0.1 的显著性水平上 $R_{OS}^2$ 显著大于 0；缩放主成分分析算法的 $R_{OS}^2$ 值为 1.68%，CW-t 检验统计量结果为 3.16，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；组合预测模型的 $R_{OS}^2$ 值为 2.94%，CW-t 检验统计量结果为 4.2，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0。

在健康行业，最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 1.21%，CW-t 检验统计量结果为 2.38，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；弹性网络算法的 $R_{OS}^2$ 值为 1.48%，CW-t 检验统计量结果为 2.6，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；主成分分析算法的 $R_{OS}^2$ 值为 0.50%，CW-t 检验统计量结果为 1.82，认为在 0.05 的显著性水平上 $R_{OS}^2$ 显著大于 0；偏最小二乘算法的 $R_{OS}^2$ 值为 1.42%，CW-t 检验统计量结果为 2.86，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；缩放主成分分析算法的 $R_{OS}^2$ 值为 1.44%，CW-t 检验统计量结果为 2.49，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；组合预测模型的 $R_{OS}^2$ 值为 1.88%，CW-t 检验统计量结果为 3.07，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0。

在其他行业，最小绝对收缩和选择算法的 $R_{OS}^2$ 值为 1.35%，CW-t 检验统计量结果为 2.31，认为在 0.05 的显著性水平上 $R_{OS}^2$ 显著大于 0；弹性网络算法的 $R_{OS}^2$ 值为 2.05%，CW-t 检验统计量结果为 2.63，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；主成分分析算法的 $R_{OS}^2$ 值为 0.13%，CW-t 检验统计量结果为 1.1，不能认为 $R_{OS}^2$ 显著大于 0；偏最小二乘算法的 $R_{OS}^2$ 值为 0.62%，CW-t 检验统计量结果为 1.9，认为在 0.05 的显著性水平上 $R_{OS}^2$ 显著大于 0；缩放主成分分析算法的 $R_{OS}^2$ 值为 1.44%，CW-t 检验统计量结果为 2.58，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0；组合预测模型的 $R_{OS}^2$ 值为 2.19%，CW-t 检验统计量结果为 2.83，认为在 0.01 的显著性水平上 $R_{OS}^2$ 显著大于 0。

对以上的结果进行总结，本文有以下三点发现：

第一，基于小波分析和机器学习方法，宏观经济指标能够有效的预测不同行业股票的超额收益。在对标普 500 指数中所包含的股票进行行业划分后，六种模型的预测效果相较基准模型均有显著提升。仅有使用主成分分析算法对“其他”行业股票进行收益率预测时，未能通过 CW-t 检验，也即主成分分析算法预测“其他”行业收益率相较取收益率均值的基准模型没有显著改善。使用其余预测模型对“消费”，“制造”，“高科技”，“健康”，“其他”行业股票收益率进行预测，均通过了 CW-t 检验，分别在不同的显著性水平上  $R_{0S}^2$  显著不为 0。在这六种预测模型当中，弹性网络算法、缩放主成分分析算法和组合算法的预测效果相较于基准模型有最大程度的提升，基于这三种预测模型，“消费”，“制造”，“高科技”，“健康”，“其他”行业的样本外 R 方值在 0.01 的显著性水平上均显著大于 0。这说明宏观经济指标不仅有对整体市场收益的预测能力，同时这种预测能力也对不同行业的股票有效。

第二，宏观经济指标对不同行业股票收益的预测能力存在异质性。对于消费行业和高科技行业，宏观指标的预测能力相对更强， $R_{0S}^2$  达到 3%-4%。这可能与消费行业/高科技行业的发展更依赖市场预期，对宏观经济的变动更加敏感。对于制造业、健康与医疗行业、其他行业的股票，宏观指标的预测能力相对较弱， $R_{0S}^2$  只有 2% 左右。对这些行业预测指标的进一步挖掘，是一个有趣的问题。

第三，不同机器学习模型在对不同行业的股票收益进行预测时，其预测能力也存在波动。比如，线性模型 Lasso 和 Elastic net 方法对于绝大多数行业有着很好的预测能力，但对于制造业的股票收益预测能力不佳。与此同时，PCA 方法对于其他行业和医疗与健康行业预测能力不佳。但组合方法通过对不同机器学习方法进行组合，减少了模型的不确定性

和预测能力的波动性，实现了偏差-方差的有效平衡，在几乎所有的行业中都实现了最佳预测结果（除高科技行业外）。因此，在对标普 500 指数根据行业进行划分的情况下，使用组合方法进行股票收益率的预测效果最佳。

#### 6.4 策略回测结果

除了择时策略，选股策略也是学术界和业界普遍关注的领域之一。与择时策略只对市场组合的持仓进行增加或减少不同，选股策略同时对不同股票的收益率进行预测，然后根据其预测值选择不同股票各自的持仓。基于不同选股策略构建的投资组合的收益情况如表 9 所示，其中已考虑 50 个基点的交易手续费，且对因变量进行了小波去噪处理。行业轮动方法中涉及到的不同行业的股票收益率基于组合预测方法进行预测。

表 9 选股策略投资组合收益情况

	return	sharpe ratio	cer	Max_Drawdown
SP500	8.8	0.57	5.2	50.4
industry rotation	14.3	0.79	9.3	32.7
industry equal_weight	9.5	0.63	6.1	47.1
industry MV	7.1	0.44	3.2	49.7
industry GMV	9.5	0.68	6.6	30.8
Cnsmr	10.1	0.66	6.6	38.2
Manuf	8.3	0.54	4.8	48.8
HiTec	9.7	0.49	3.9	77.8
Hlth	10.3	0.64	6.5	35.3
Other	8.8	0.49	3.9	68.3

使用行业轮动方法构建投资组合，得到年化收益率为 14.26%，得到夏普比率为 0.79，得到年化确定性等价收益为 9.33，得到最大回撤为 32.75%；同时，等权重策略，也被称为 1/N 投资策略，也是被广泛使用的一种投资策略。DeMiguel, Garlappi & Uppal (2007) 调查了 14 种投资组合构建方法，发现没有任何一种能够在所有数据集上击败等

权重策略。使用行业等权重方法构建投资组合，投资组合的年化收益率为 9.46%，夏普比率为 0.63，年化确定性等价收益为 6.09，最大回撤为 47.12%；本文也将市场组合作为一种基准模型，其年化收益率为 8.81%，夏普比率为 0.57，年化确定性等价收益为 5.18，最大回撤为 50.4%；使用均值-方差有效前沿方法构造的投资组合，年化收益率为 7.09%，夏普比率为 0.44，年化确定性等价收益为 3.23，最大回撤为 49.71%；而最小方差组合的年化收益率为 9.5%，夏普比率为 0.68，年化确定性等价收益为 6.56，最大回撤为 30.76%；有许多基金专注于特定行业的投资，因此本文也比较了单个行业的投资收益。消费行业组合的年化收益率为 10.09%，夏普比率为 0.66，年化确定性等价收益为 6.6%，最大回撤为 38.19%；制造业的年化收益率为 8.35%，夏普比率为 0.54，年化确定性等价收益为 4.8，最大回撤为 48.77%；高科技行业组合的年化收益率为 9.71%，夏普比率为 0.49，年化确定性等价收益为 3.86，最大回撤为 77.8%，高科技行业的几次泡沫破灭极大地增加了波动率，从而降低了其投资收益；对于医疗与健康行业，年化收益率为 10.29%，夏普比率为 0.64，年化确定性等价收益为 6.45，最大回撤为 35.35%；最后是其其他行业组合，年化收益率为 8.84%，夏普比率为 0.49，年化确定性等价收益为 3.92，最大回撤为 68.27%。

在所有构建投资组合的方法中，行业轮动方法获得了最高的年化收益率、年化确定性等价收益、年化夏普比率与几乎最低的最大回撤。这是因为行业轮动方法利用了行业的轮动特性，即不同行业在不同时间段内的表现会发生变化。行业轮动方法通过对行业表现的监测和分析，能够及时调整投资组合的配置，使投资组合在市场的不同阶段获得更好的收益。相比于长期持有某些行业的股票，如使用单一行业持有的方法构建投资组合，行业轮

动方法能够更加灵活地应对市场的变化，减少投资组合的风险。与均值-方差有效前沿方法和最小方差组合方法相比，行业轮动投资组合的构建过程中，考虑到了宏观经济对同行业股票市场的异质性影响，某种意义上实现了择时和选股的结合。此外，行业轮动方法能够减少投资组合的相关性，从而实现更好的分散化。这种方法通过投资不同行业的股票，能够减少投资组合内股票之间的相关性，从而降低投资组合的波动性

## 6.5 累积财富积累曲线

图 5 选股策略的收益率曲线

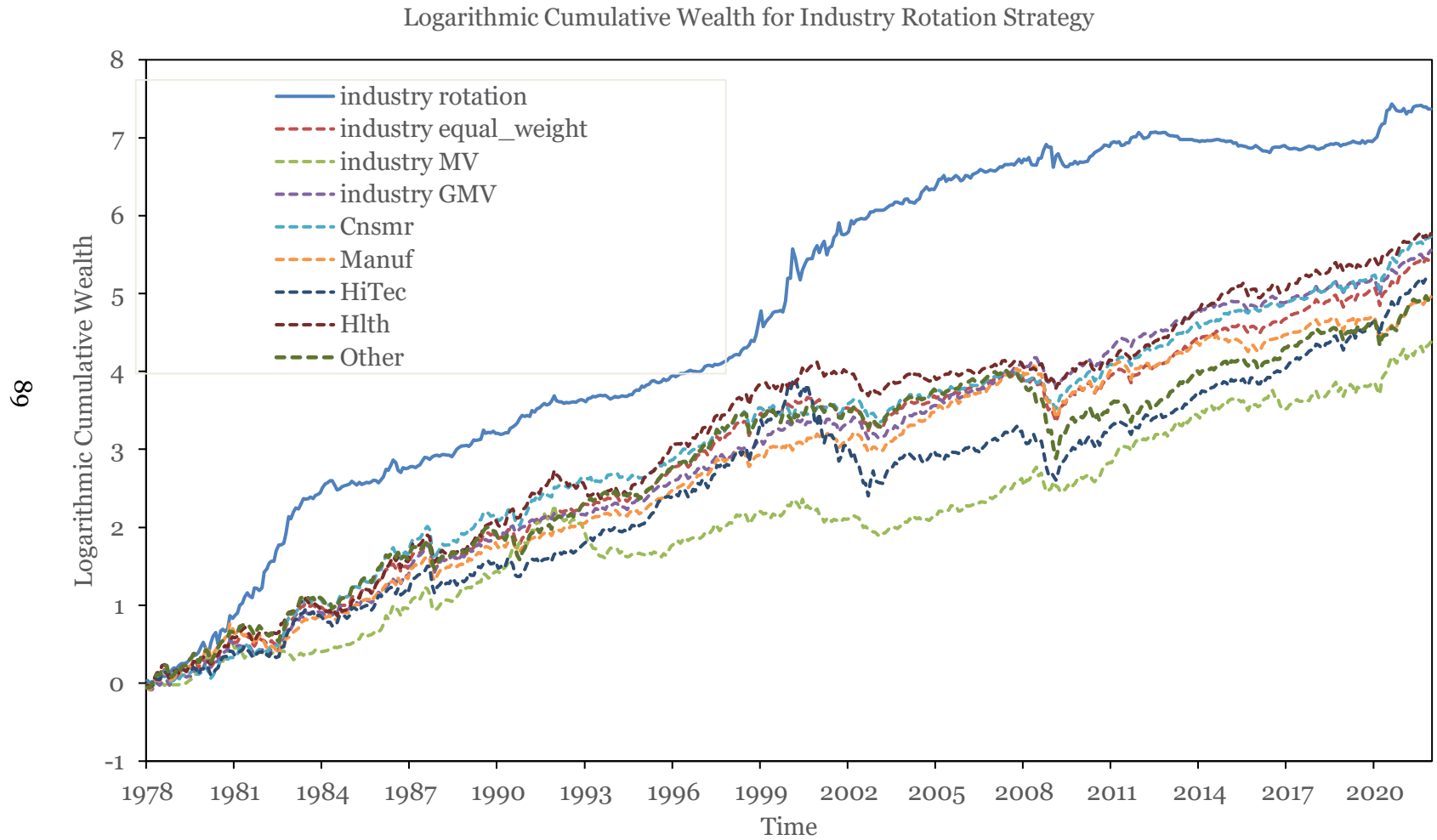


图 5 为基于行业的不同方法构建投资组合得到的收益率对数曲线。基于行业轮动方法构建的投资组合（蓝色实线）相较于其他投资组合拥有显著更高的累计收益率曲线。“制造”、“行业等权”、“健康”等方法构建的投资组合（橙色虚线、深色虚线、棕色虚线）之间的收益差异并不是十分显著。在标普 500 指数发生下跌时，这些投资组合的收益率同步发生下跌，没有实现对风险的有效对冲。

观察除行业轮动方法外其余各个投资组合得到的累计收益率曲线，可以发现收益率曲线在 1988 年、1991 年、2002 年和 2008 年附近均出现明显的大幅下跌。回顾这些时间段发生的历史事件，可以找到这些年份附近发生大幅下跌的原因：1987 年 10 月 19 日，美国股市发生了历史上最大的单日跌幅，被称为“黑色星期一”，标普 500 指数当日下跌 22.6%。此后的一年内，美国经济进入了一段调整期，特别是房地产和储蓄贷款危机的影响开始显现。由于市场对于经济衰退的担忧和恐慌情绪，标普 500 指数在 1988 年下跌了 11.4%。1990 年 8 月，伊拉克入侵科威特，引发了波斯湾战争。此后，国际原油价格大幅上涨，导致通货膨胀率上升。此外，美国经济的增长也放缓了。由于市场对于经济增长的担忧，标普 500 指数在 1991 年下跌了 6.6%。2001 年 9 月 11 日，恐怖分子在美国发动了针对世界贸易中心和五角大楼的袭击。这次袭击引发了全球经济的恐慌，导致许多国家的股市下跌。此外，互联网泡沫破灭，也导致了許多科技公司的股价下跌。由于这些因素的影响，标普 500 指数在 2002 年下跌了 22.1%。2007 年底，美国次贷危机爆发，许多银行和金融机构因为持有大量次贷债券而陷入了危机。次贷危机引发了金融市场的恐慌，导致全球股市下跌。此外，全球经济增长放缓，加上高油价和粮食价格上涨，导致通货膨胀率上升。由于这些因素的影响，标普 500 指数在 2008 年下跌了 38.5%。

可以发现使用“行业轮动”方式构建的投资组合通过构建市场风险对冲机制，在这些标普 500 指数发生大幅下跌的时间点，仍然保持着较为平稳的累计收益率曲线。同时，该方法充分利用了被其他模型所忽略的宏观因子信息，对不同行业的收益率进行了准确的预测。该方法保证了投资组合的稳定性和长期盈利能力，在不断变化的市场中帮助投资者更快地应对市场变化和风险，并且在不同的市场条件下实现良好的收益。

## 七、预测框架与方法的优缺点分析

本文通过科学的方法论和实证分析，揭示了大数据和机器学习在经济预测中的巨大潜力。然而，笔者也清醒地认识到，虽然这些新的技术和方法提供了前所未有的视角和工具，但是它们也存在自身的挑战和局限性，在实际的应用中，需要更为细心和谨慎地考虑和处理。

首先，总结本文提出的预测框架和方法的优点：

第一，本文的预测模型采用了海量数据进行训练和预测，符合当前大数据时代的基本要求。在当下所处的这个信息爆炸的时代，无论是结构化还是非结构化的数据都呈现出指数级的增长。在这种背景下，本文的预测模型能够有效地利用这些海量数据，揭示并利用其中隐藏的、复杂的模式，更全面、更准确地反映市场和经济活动的真实状况。这一点无疑是传统的、依赖于少量数据的预测模型所无法比拟的。大数据不仅可以帮助捕捉到更细致、更复杂的市场动态，还可以通过数据的深度和广度，帮助对市场进行更精细的划分，提高预测模型的准确性和鲁棒性。因此本文的方法充分利用了大数据的优势，提供了更为深入和精确的预测结果。

第二，本文的预测方法采用了组合机器学习模型。在实际的研究和预测中，单一模型常常因为其固有的先验条件和模型设定的限制，存在模型误设的问题。而且，单一模型的预测能力常常是时变的，会受到预测期间经济环境、市场状态等因素的影响。因此，本文选择了组合机器学习模型的方式，对不同的机器学习模型进行了加权。这样做的优点在于，可以充分利用模型的多样性，减少模型误设的风险，同时也可以提高预测的稳定性和精度。

第三，本文应对了股票收益率数据的挑战。月度股票收益率的数据往往存在剧烈波动和大量噪音，这给本文的研究带来了很大的挑战。为了解决这个问题，本文采用了小波频域分解技术，对股票收益率数据进行了频域分解，以提取不同频率和时间尺度上的成分。这一方法有助于从复杂、嘈杂的数据中筛选出真正有价值的信息，从而提高了预测的准确性。

第四，本文的预测方法有着广泛的应用性。本文不仅将预测模型应用在了市场收益预测上，实现了一种有效的指数增强策略，同时，也将预测模型应用于了不同行业的股票收益预测。通过预测各行业的股票收益，本文构建了买入-卖出投资组合，实现了横截面上的选股套利，这是一个相对于其他文献的创新点。

当然，股票收益预测是非常困难的，也必须面对一些不可避免的挑战和局限性。以下是本文提出的预测框架和方法的一些局限性：

第一，本文采用的海量数据主要来源于已经发表的公开数据集，这无疑在一定程度上限制了预测模型对新的和实时的经济活动的理解和预测。随着论文的出版和时间的推移，这些宏观经济指标的预测能力可能会逐渐下降，这是一个需要深入关注的问题。针对这个问题，笔者在未来的研究中，可以考虑采用更高频率的宏观数据，或者引入更多的非结构化数据，如文本信息、图像信息等。这些数据资源不仅可以提供更多样的信息，也可能会帮助预测模型更好地捕捉市场和经济的动态变化，提高预测精度。

第二，本文的预测模型使用了递归扩展窗口的方式来估计参数。这在某种程度上忽视了在金融数据中常常出现的一种复杂现象，即结构变迁。结构变迁就是指模型中的因变量和自变量之间的关系随着时间的推移而发生变化。这是金融预测研究一个棘手的问题，它

带来的挑战超越了本文的讨论范围。但是，这并不意味着可以忽视其存在。相反，在未来的研究中，笔者需要进一步深化对这一问题的理解，并在模型设计和参数估计中，尝试寻找能够有效应对这一问题的解决策略。这样，预测模型在面对市场和经济的动态变化时，才能保持较高的稳定性和预测精度。

尽管本文的预测框架和方法存在这些挑战和局限性，但笔者相信，随着研究的深入和新技术、新方法的发展，将能够更好地解决这些问题，进一步提高预测的准确性和鲁棒性。

## 参考文献

1. Bansal, Ravi, and Amir Yaron. "Risks for the long run: A potential resolution of asset pricing puzzles." *The journal of Finance* 59, no. 4 (2004): 1481-1509.
2. Welch, Ivo, and Amit Goyal. "A comprehensive look at the empirical performance of equity premium prediction." *The Review of Financial Studies* 21, no. 4 (2008): 1455-1508.
3. Driesprong, Gerben, Ben Jacobsen, and Benjamin Maat. "Striking oil: another puzzle?." *Journal of financial economics* 89, no. 2 (2008): 307-327.
4. Bollerslev, Tim, George Tauchen, and Hao Zhou. "Expected stock returns and variance risk premia." *The Review of Financial Studies* 22, no. 11 (2009): 4463-4492.
5. Brogaard, Jonathan, Andrew Detzel, and Phong TH Ngo. "Inequality and risk premia." Available at SSRN (2015).
6. Huang, Dashan, Fuwei Jiang, Jun Tu, and Guofu Zhou. "Investor sentiment aligned: A powerful predictor of stock returns." *The Review of Financial Studies* 28, no. 3 (2015): 791-837.
7. Rapach, David E., Matthew C. Ringgenberg, and Guofu Zhou. "Short interest and aggregate stock returns." *Journal of Financial Economics* 121, no. 1 (2016): 46-65.
8. Manela, Asaf, and Alan Moreira. "News implied volatility and disaster concerns." *Journal of Financial Economics* 123, no. 1 (2017): 137-162.
9. Jiang, Fuwei, Joshua Lee, Xiumin Martin, and Guofu Zhou. "Manager sentiment and stock returns." *Journal of Financial Economics* 132, no. 1 (2019): 126-149.
10. Wang, Yudong, Zhiyuan Pan, Chongfeng Wu, and Wenfeng Wu. "Industry equity correlation: A powerful predictor of stock returns." *Journal of Empirical Finance* 59 (2020): 1-24.
11. Chen, Lin, Zhilin Qiao, Minggang Wang, Chao Wang, Ruijin Du, and Harry Eugene Stanley. "Which artificial intelligence algorithm better predicts the Chinese stock market?." *IEEE Access* 6 (2018): 48625-48633.
12. Zhang, Kang, Guoqiang Zhong, Junyu Dong, Shengke Wang, and Yong Wang. "Stock market prediction based on generative adversarial network." *Procedia computer science* 147 (2019): 400-406.

13. Huang, Yuxuan, Luiz Fernando Capretz, and Danny Ho. "Machine learning for stock prediction based on fundamental analysis." In 2021 IEEE Symposium Series on Computational Intelligence (SSCI), pp. 01-10. IEEE, 2021.
14. Cochrane, John H. "The dog that did not bark: A defense of return predictability." *The Review of Financial Studies* 21, no. 4 (2008): 1533-1575.
15. Rapach, David, and Guofu Zhou. "Forecasting stock returns." In *Handbook of economic forecasting*, vol. 2, pp. 328-383. Elsevier, 2013.
16. Huang, Dashan, and Guofu Zhou. "Upper bounds on return predictability." *Journal of Financial and Quantitative Analysis* 52, no. 2 (2017): 401-425.
17. Campbell, John Y. "Asset pricing at the millennium." *The Journal of Finance* 55, no. 4 (2000): 1515-1567.
18. Chen, Shiu-Sheng. "Predicting the bear stock market: Macroeconomic variables as leading indicators." *Journal of Banking & Finance* 33, no. 2 (2009): 211-223.
19. Nti, Kofi O., Adebayo Adekoya, and Benjamin Weyori. "Random forest based feature selection of macroeconomic variables for stock market prediction." *American Journal of Applied Sciences* 16, no. 7 (2019): 200-212.
20. Chowdhury, Shah, Abu Taher Mollik, and M. Selim Akhter. "Does predicted macroeconomic volatility influence stock market volatility? Evidence from the Bangladesh capital market." PhD diss., University of South Australia, 2006.
21. Ludvigson, Sydney C., and Serena Ng. "The empirical risk–return relation: A factor analysis approach." *Journal of financial economics* 83, no. 1 (2007): 171-222.
22. Bialkowski, Jędrzej, Katrin Gottschalk, and Tomasz Piotr Wisniewski. "Political orientation of government and stock market returns." *Applied Financial Economics Letters* 3, no. 4 (2007): 269-273.
23. Jacobsen, Ben, Ben R. Marshall, and Nuttawat Visaltanachoti. "Stock market predictability and industrial metal returns." *Management Science* 65, no. 7 (2019): 3026-3042.
24. Salisu, Afees A., and Xuan Vinh Vo. "Predicting stock returns in the presence of COVID-19 pandemic: The role of health news." *International Review of Financial Analysis* 71 (2020): 101546.
25. Han, Yufeng, Ai He, David Rapach, and Guofu Zhou. "Firm characteristics and expected stock returns." Available at SSRN (2019).

26. Wold, Herman. "Estimation of principal components and related models by iterative least squares." *Multivariate analysis* (1966): 391-420.
27. Wold, Herman. "Path models with latent variables: The NIPALS approach." In *Quantitative sociology*, pp. 307-357. Academic Press, 1975.
28. Kelly, Bryan, and Seth Pruitt. "Market expectations in the cross-section of present values." *The Journal of Finance* 68, no. 5 (2013): 1721-1756.
29. Kelly, Bryan, and Seth Pruitt. "The three-pass regression filter: A new approach to forecasting using many predictors." *Journal of Econometrics* 186, no. 2 (2015): 294-316.
30. Rapach, David E., Jack K. Strauss, and Guofu Zhou. "International stock return predictability: What is the role of the United States?." *The Journal of Finance* 68, no. 4 (2013): 1633-1662.
31. Gu, Shihao, Bryan T. Kelly, and Dacheng Xiu. "Empirical asset pricing via machine learning." *Chicago Booth Research Paper* 18-04 (2019): 2018-09.
32. Chincó, Alex, Adam D. Clark-Joseph, and Mao Ye. "Sparse signals in the cross-section of returns." *The Journal of Finance* 74, no. 1 (2019): 449-492.
33. Freyberger, Joachim. "On completeness and consistency in nonparametric instrumental variable models." *Econometrica* 85, no. 5 (2017): 1629-1644.
34. Kozak, Serhiy, Stefan Nagel, and Shrihari Santosh. "Shrinking the cross-section." *Journal of Financial Economics* 135, no. 2 (2020): 271-292.
35. Kong, Aiguo, David Rapach, Jack Strauss, Jun Tu, and Guofu Zhou. "How predictable are components of the aggregate market portfolio." *Research Collection Lee Kong Chian School of Business* (2009).
36. Green, Jeremiah, John RM Hand, and X. Frank Zhang. "The characteristics that provide independent information about average US monthly stock returns." *The Review of Financial Studies* 30, no. 12 (2017): 4389-4436.
37. Hong, Yongmiao, Fuwei Jiang, Lingchao Meng, and Bowen Xue. "Forecasting Inflation with Economic Narratives and Machine Learning." Available at SSRN 4175749 (2022).
38. Markowitz, Harry. "The utility of wealth." *Journal of political Economy* 60, no. 2 (1952): 151-158.

39. Barry, Christopher B. "Portfolio analysis under uncertain means, variances, and covariances." *The Journal of Finance* 29, no. 2 (1974): 515-522.
40. Bawa, Vijay S., Stephen J. Brown, and Roger W. Klein. "Estimation risk and optimal portfolio choice." (1979).
41. Jobson, J. D. "Improved estimation for Markowitz portfolios using James-Stein type estimators." In *Proceedings of the American Statistical Association, Business and Economics Statistics Section*, vol. 71, pp. 279-284. 1979.
42. Jobson, J. David, and Bob Korkie. "Estimation for Markowitz efficient portfolios." *Journal of the American Statistical Association* 75, no. 371 (1980): 544-554.
43. Detemple, Jerome, and Philippe Jorion. "Option listing and stock returns: An empirical analysis." *Journal of Banking & Finance* 14, no. 4 (1990): 781-801.
44. Pástor, Luboš, and Robert F. Stambaugh. "Comparing asset pricing models: an investment perspective." *Journal of Financial Economics* 56, no. 3 (2000): 335-381.
45. Pastor, Lubos, and Robert F. Stambaugh. "Evaluating and investing in equity mutual funds." (2000).
46. Goldfarb, Donald, and Garud Iyengar. "Robust portfolio selection problems." *Mathematics of operations research* 28, no. 1 (2003): 1-38.
47. Garlappi, Lorenzo, Raman Uppal, and Tan Wang. "Portfolio selection with parameter and model uncertainty: A multi-prior approach." *The Review of Financial Studies* 20, no. 1 (2007): 41-81.
48. Kan, Raymond, and Guofu Zhou. "Optimal portfolio choice with parameter uncertainty." *Journal of Financial and Quantitative Analysis* 42, no. 3 (2007): 621-656.
49. Craig MacKinlay, A., and Luboš Pástor. "Asset pricing models: Implications for expected returns and portfolio selection." *The Review of financial studies* 13, no. 4 (2000): 883-916.
50. Best, Michael J., and Robert R. Grauer. "Positively weighted minimum-variance portfolios and the structure of asset expected returns." *Journal of Financial and Quantitative Analysis* 27, no. 4 (1992): 513-537.
51. Chan, Louis KC, Jason Karceski, and Josef Lakonishok. "On portfolio optimization: Forecasting covariances and choosing the risk model." *The review of Financial studies* 12, no. 5 (1999): 937-974.

52. Ledoit, Olivier, and Michael Wolf. "A well-conditioned estimator for large-dimensional covariance matrices." *Journal of multivariate analysis* 88, no. 2 (2004): 365-411.
53. Frost, Peter A., and James E. Savarino. "For better performance: Constrain portfolio weights." *Journal of portfolio management* 15, no. 1 (1988): 29.
54. Chopra, Vijay K., and William T. Ziemba. "The effect of errors in means, variances, and covariances on optimal portfolio choice." In *Handbook of the fundamentals of financial decision making: Part I*, pp. 365-373. 2013.
55. Jagannathan, Ravi, and Tongshu Ma. "Risk reduction in large portfolios: Why imposing the wrong constraints helps." *The journal of finance* 58, no. 4 (2003): 1651-1683.
56. DeMiguel, Victor, Lorenzo Garlappi, and Raman Uppal. "Optimal versus naive diversification: How inefficient is the 1/N portfolio strategy?." *The review of Financial studies* 22, no. 5 (2009): 1915-1953.
57. Neely, Christopher J., David E. Rapach, Jun Tu, and Guofu Zhou. "Forecasting the equity risk premium: the role of technical indicators." *Management science* 60, no. 7 (2014): 1772-1791.
58. Clark, Todd E., and Kenneth D. West. "Approximately normal tests for equal predictive accuracy in nested models." *Journal of econometrics* 138, no. 1 (2007): 291-311.
59. Hong, Yongmiao, Fuwei Jiang, Lingchao Meng, and Bowen Xue. "Forecasting Inflation with Economic Narratives and Machine Learning." Available at SSRN 4175749 (2022).
60. Baker, Malcolm, and Jeffrey Wurgler. "Investor sentiment in the stock market." *Journal of economic perspectives* 21, no. 2 (2007): 129-151.
61. Kang, Jie, Fuwei Jiang, and Zhifeng Dai. "Stock Return Predictability in Frequency Domain." Available at SSRN 4266050.
62. Clark, Todd E., and Kenneth D. West. "Using out-of-sample mean squared prediction errors to test the martingale difference hypothesis." *Journal of econometrics* 135, no. 1-2 (2006): 155-186.

63. Rapach, David E., Jack K. Strauss, and Guofu Zhou. "Out-of-sample equity premium prediction: Combination forecasts and links to the real economy." *The Review of Financial Studies* 23, no. 2 (2010): 821-862.
64. Rapach, David E., and Guofu Zhou. "Time-series and cross-sectional stock return forecasting: New machine learning methods." *Machine learning for asset management: New developments and financial applications* (2020): 1-33.
65. Granger, Clive WJ. "Investigating causal relations by econometric models and cross-spectral methods." *Econometrica: journal of the Econometric Society* (1969): 424-438.
66. Wold, Svante, Michael Sjöström, and Lennart Eriksson. "PLS-regression: a basic tool of chemometrics." *Chemometrics and intelligent laboratory systems* 58, no. 2 (2001): 109-130.
67. Huang, Dashan, Fuwei Jiang, Kunpeng Li, Guoshi Tong, and Guofu Zhou. 2022. "Scaled PCA: A new approach to dimension reduction" *Management Science* 68, no. 3 (2022): 1678-1695.
68. Chen J, Tang G, Yao J, Zhou G. "Investor Attention and Stock Returns" *Journal of Financial and Quantitative Analysis* 57, no. 2 (2022): 455-484.
69. Baker, Malcolm, and Jeffrey Wurgler. "Investor sentiment and the cross-section of stock returns." *The journal of Finance* 61, no. 4 (2006): 1645-1680.
70. McCracken, Michael W., and Serena Ng. "FRED-MD: A monthly database for macroeconomic research." *Journal of Business & Economic Statistics* 34, no. 4 (2016): 574-589.
71. Ludvigson, Sydney C., and Serena Ng. "Macro factors in bond risk premia." *The Review of Financial Studies* 22, no. 12 (2009): 5027-5067.