

CITY UNIVERSITY OF HONG KONG
香港城市大學

**Portfolio Investment Strategy Based on
Quantitative Industry Selection and Machine
Learning Methods**

**基於量化產業選擇和機器學習方法的投資
組合策略研究**

Submitted to
College of Business
商學院
in Partial Fulfillment of the Requirements
for the Degree of Doctor of Business Administration
工商管理博士學位

by

Chen Dan
陳丹

March 2022
二零二二年三月

摘要

中國經濟及證券市場受政策影響非常大，近幾十年來中國經濟的高速發展，在一定程度上也得益於國家正確的產業政策。然而，作為經濟發展晴雨錶的證券市場的發展卻並不如人意，沒有與經濟同步發展，通過證券投資並沒有分享到經濟高速發展所帶著的收益，其中很大一部分原因就是投資只是追熱點證券與板塊，而忽視了對相關行業進行深入研究。行業發展有一定的趨勢性和規律性，行業研究是證券投資的手段，具有前瞻性與穩定性，尤其是金融機構的行業研究對相關產業政策與行業發展會有很好的把握。本文從行業研究的角度研究中國證券市場的投資問題。現有的行業研究基本上是基於宏觀經濟與產業政策的定性研究，缺乏定量分析，本文從量化的角度研究行業投資問題。

本文首先對中國資本市場及行業發展的現狀進行回顧與分析；之後研究中國證券市場基本情況及行業投資的現狀進行深入剖析；然後對國內外相關研究進行梳理與分析。在此基礎上主要進行了以下四部分工作。

第一，在借鑒綜合學者行業研究的基礎上，運用數量分析方法，構建行業發展指數。行業發展指數綜合了行業內在資訊及相關產業的社會關注程度。包括行業的基本面資訊、技術指標，以及金融機構研究關注程度等等，能夠反應行業的現狀及發展趨勢。然後運用行業指數選擇投資的領域。

第二，接下來根據選擇的投資行業選擇相應的股票，構建投資組合。對特定行業運用機器學習方法，運用行業內股票的基本面與技術面資訊構建投資組合。投資組合具有較強的代表性，能夠反應行業的發展趨勢。選定之後對投資組合就行優化。本文多種機器學習方法選擇行業內證券，在比較各個

模型的優劣基礎上，提出集成機器學習方法就行業內股票選擇，對擇股能力進一步提升。研究行業約束條件下的投資組合優化問題。

第三，研究投資的風險度量與控制問題。本文拓展了業界普遍採用的風險度量指標最大回撤率，將最大回撤賦予了概率度量，創新性的提出了 DaR 指標，採用模擬的方法產生最大回撤曲面，並給出實際風險管理與控制操作方法。

第四，給出投資策略的實際運行方法。在前面根據行業發展指數的基礎上，運用機器學習方法選擇行業內股票，然後根據多種技術面資訊進行交易時機的選擇，接下來進行投資組合優化及風險控制。

最後進行模擬交易，無論是行業發展指數，還是運用機器學習方法選擇行業內股票對投資都有顯著的提升，都會增加收益，減少風險。

創新點：本文在以下幾個方面進行了創新。首先，構建綜合行業內外因素的行業發展指數，根據行業發展指數選擇投資行業；然後，針對特定行業運用機器學習方法選擇股票構建投資組合，反應行業發展趨勢；第三，提出改進的風險度量指標 DaR，給出新指標的定義及計算方法。最後，給出投資策略的實踐方法。

本文在解決行業發展評估、投資組合構建、風險管控等幾個關鍵問題的基礎上，構建較為全面的行業量化投資框架，對行業研究與投資實務具有較強的實際價值；同時在行業發展量化研究、以及差異化風險管控方面，也具有一定的理論與方法上的創新。

關鍵字：行業發展指數；機器學習；風險度量；量化投資；Stacking 集成學習

ABSTRACT

China's economy and securities market are greatly affected by policies. The rapid development of China's economy in recent decades also benefits from the correct national industrial policies to a certain extent. However, as a barometer of economic development, the development of the securities market is not satisfactory. It does not develop synchronously with the economy. Through securities investment, it does not share the benefits brought by the rapid economic development. A large part of the reason is that investment only pursues the hot securities and plates, and ignores the in-depth research on relevant industries. Industry development has certain trends and regularity. Industry research is a means of securities investment, which is forward-looking and stable. Especially the industry research of financial institutions will have a good grasp of relevant industrial policies and industry development. This paper studies the investment in China's securities market from the perspective of industry research. The existing industry research is basically based on the qualitative research of macroeconomic and industrial policy, lack of quantitative analysis. This dissertation studies the industry investment from the perspective of quantification.

Firstly, this paper reviews and analyzes the current situation of China's capital market and industry development; Then it studies the basic situation of China's securities market and the current situation of industry investment; Then it combs and analyzes the relevant research at home and abroad. On this basis, the following four parts are mainly carried out.

First, based on the industry research of comprehensive scholars, the quantitative analysis method is used to construct the industry development index. The industry development index integrates the internal information of the industry and the social attention of related industries. Including the basic information of the industry, technical indicators, and the research attention of financial institutions, which can reflect the current situation and development trend of the industry. Then use the industry index to select the investment field.

Second, select the corresponding stocks according to the selected investment industry to build the investment portfolio. Use machine learning methods for specific industries, and use the fundamental and technical information of stocks in the industry to build a portfolio. The portfolio is highly representative and can reflect the development trend of the industry. After selection, the portfolio will be optimized. In this paper, a variety of machine learning methods are used to select securities in the industry. Based on the comparison of the disadvantages and advantages of each model, an integrated machine learning method is proposed to further improve the stock selection ability in the industry. The problem of portfolio optimization under industry constraints is studied.

Third, research the risk measurement and control of investment. This paper expands the risk measurement index maximum pullback rate commonly used in the industry, endows the maximum pullback with probability measurement, innovatively puts forward Dar index, uses simulation method to generate the maximum pullback surface, and gives the actual risk management and control operation method.

Fourth, the actual operation method of investment strategy is given. Based on the previous industry development index, the machine learning method is used to select the stocks in the industry, and then the trading timing is selected according to a variety of technical information, followed by portfolio optimization and risk control.

Finally, simulated trading, whether industry development index or using machine learning method to select stocks in the industry, will significantly improve the investment, increase the income and reduce the risk.

Innovation: This dissertation makes innovations in the following aspects. Firstly, build an industry development index integrating internal and external factors of the industry, and select the investment industry according to the industry development index; Then, according to the specific industry, the machine learning method is used to select stocks to build a portfolio to reflect the development trend of the industry; Thirdly, an improved risk measurement index Dar is proposed, and the definition and calculation method of the new index are given. Finally, the

practical method of investment strategy is given.

On the basis of solving several key problems such as industry development evaluation, portfolio construction and risk control, this dissertation constructs a more comprehensive industry quantitative investment framework, which has strong practical value for Industry Research and investment practice; At the same time, it also has some theoretical and methodological innovations in quantitative research on industry development and differentiated risk management and control.

Key words: Industry development index; Machine learning; Risk measurement; Quantitative investment; Stacking ensemble learning

致謝

感恩！讀博是壹個漫長的群體實踐，為了壹個不確定的研究，導師、同學、家人、還有很多愛我的親人和朋友付出了大量的關愛！感謝我的導師馬躍教授、芮明傑教授給予我的指導和包容！感謝我的助研郭進老師和同門譚凌波老師的傾力幫助！感謝 2015 級全體同學這個溫暖的集體和領頭大哥孫屹崕班長！感謝 Angel CHAN 老師和白碧湖老師的熱情幫助！最後，感謝父母、家人和親人默默的等待和期盼！

讀博也是壹件悲喜交集的事情，探尋壹個成為我自己感興趣能解決的問題就是壹個問題！但是壹旦條件成熟，小概率事件也成為敲開大門的那塊磚。讀博初衷僅僅是為人生打開壹扇門，但是讀博期間也打開了很多扇窗，擁抱不確定性的未來，篤定前行。

目錄

摘要	i
ABSTRACT	iii
Qualifying Panel and Examination Panel	vi
致謝	vii
圖表目錄	x
第一章 引言	1
1.1 從中國證券市場的角度.....	1
1.2 從我國的產業政策的角度.....	4
1.3 從技術發展的角度.....	4
1.4 研究意義.....	5
第二章 文獻綜述	6
2.1 經典理論及其挑戰.....	6
2.1.1 經典投資模型與有效市場理論	6
2.1.2 動量投資效應和反轉投資效應對傳統理論的挑戰	8
2.2 目前的解釋及我們的視角	11
2.3 政策在中國證券市場中是非常重要的影響因素	13
2.4 量化投資在中國市場的應用	14
第三章 行業發展指數研究	16
3.1 行業內在因素	16
3.1.1 估值因素	16
3.1.2 行業動量	17
3.1.3 市場活躍程度	18
3.1.4 市場回報	18
3.2 行業關注度	20
3.3 行業發展指數	30
第四章 行業內機器學習選擇股票	58
4.1 指標選擇	58
4.2 機器學習模型預測行業內優質股票	59
4.2.1 隨機森林	59
4.2.2 AdaBoost	61
4.3 運用集成學習提高選股能力	65
4.3.1 不同模型捕捉善於不同方面的資訊	65
4.3.2 Stacking 集成	71
第五章 風險度量與控制	76
5.1 最大回撤率	76
5.2 最大回撤模擬	77
5.3 不同置信水準的回撤（DaR）	82
5.4 不同波動性的回撤（DaR）	87
5.5 DaR 的投資應用	91
第六章 投資應用及效果	95
6.1 交易時機的選擇	95
6.2 投資組合優化	100

6.3 投資效果模擬	102
第七章 總結與展望	105
7.1 總結	105
7.2 問題及展望	106
參考文獻	108
附錄一 關鍵程式代碼	116
附錄二 投資組合優化	127

圖表目錄

圖 1 中國市場重要行業指數	1
圖 2 滬深 300 指數與食品飲料行業指數	2
圖 3 滬深 300 指數與交通運輸行業指數	2
圖 4 主要行業價格變化	3
圖 5 主要行業月度價格變化	4
圖 6 2021 年醫藥行業市盈率	16
圖 7 2021 年醫藥行業市淨率	17
圖 8 2021 年醫藥行業上漲比率	18
圖 9 醫藥行業的換手率	18
圖 10 2021 年醫藥行業收益率	19
圖 11 網路資訊資源	21
圖 12 2021 年全年行業熱詞雲圖	23
圖 13 2021 年全年行業詞頻統計	23
圖 14 2021 年 10 月行業熱詞雲圖	24
圖 15 2021 年 10 月行業詞頻統計	25
圖 16 2021 年 11 月行業熱詞雲圖	25
圖 17 2021 年 11 月行業詞頻統計	26
圖 18 電腦行業研究報告頻率圖	26
圖 19 電腦行業研究報告標準頻率圖	27
圖 20 銀行業研究報告標準頻率圖	27
圖 21 醫藥行業研究報告標準頻率圖	28
圖 22 汽車行業研究報告標準頻率圖	28
圖 23 電腦行業研究報告標準頻率圖	29
圖 24 鋼鐵行業研究報告標準頻率圖	29
圖 25 典型行業研究報告標準頻率圖	30
圖 26 醫藥行業發展指數	32
圖 27 醫藥行業價格指數	32
圖 28 隨機森林模型構造流程圖	60
圖 29 AdaBoost 基本結構流程圖	62
圖 30 各個因數在隨機森林模型中的重要性（從大到小排序）	65
圖 31 各個因數在 GBDT 模型中的重要性（從大到小排序）	67
圖 32 各個因數在 XGB 模型中的重要性（從大到小排序）	68
圖 33 各個因數在 AdaBoost 模型中的重要性（從大到小排序）	70
圖 34 多種機器學習方法集成技術	72
圖 35 神經網路結構圖	73
圖 36 relu 曲線圖	73
圖 37 Stacking 的 ROC 曲線	75
圖 38 投資的淨值與最大回撤	76
圖 39 DaR 示意圖	77
圖 40 模擬 1000 次最大回撤	78
圖 41 模擬 1000 次最大回撤均值	78
圖 42 模擬 10000 次最大回撤	79

圖 43 最大回撤的分佈	80
圖 44 最大回撤的均值	80
圖 45 50*50000 次模擬	81
圖 46 50*50000 次模擬的最大回撤的均值	82
圖 47 標準差 20%最大回撤的分佈	83
圖 48 標準差 50%最大回撤的分佈	83
圖 49 標準差 10%最大回撤的分佈	84
圖 50 最大回撤曲面	85
圖 51 不同標準差與最大回撤關係圖	85
圖 52 不同置信水準的標準差與最大回撤關係圖	86
圖 53 標準差 0.02，對應不同置信區間的回撤	87
圖 54 不同波動率下的 DaR	88
圖 55 標準差最大回撤聯合概率分佈圖	89
圖 56 不同標準差對應的最大回撤的分佈曲線	90
圖 57 不同標準差對應的分佈曲線對比圖	91
圖 58 根據 DaR 尋找波動率	92
圖 59 根據行業發展指數投資	103
圖 60 行業發展指數加機器學習選股投資	104
 表 1 醫藥行業相關數據	19
表 2 醫藥行業基本數據統計特徵	20
表 3 大資料捕捉金融機構研究報告關鍵內容	21
表 4 典型研究報告標準頻率統計特徵表	29
表 5 選股指標	58
表 6 混淆矩陣	63
表 7 機器學習效果關鍵評價指標對比	64
表 8 各個因數在隨機森林模型中的重要性	66
表 9 各個因數在 GBDT 模型中的重要性	67
表 10 各個因數在 XGB 模型中的重要性	69
表 11 各個因數在 AdaBoost 模型中的重要性	70
表 12 Stacking 集成演算法與基礎模型比較	74
表 13 最大回撤與標準差	86
表 14 不同波動率下的置信區間的最大回撤（DaR）	88