

**Final Programming Exam. (part 1)**  
**Deep Learning, 2022**

姓名：吳佳展      學號：110318517

● **Project Title：**

**結合趨勢的深度強化學習股票交易策略**

● **Introduction：**

使用深度學習預測股票價格或趨勢運動有許多研究成果。BAO 等[2]將小波變換、堆線式自編碼器(SAEs)和 LSTM 相互結合，首先對股票價格時間序列進行小波變換分解以消除噪點，接著應用 SAEs 生成深層高階特徵，輸入至 LSTM 中預測第二天的收盤價。Cai 等[3]提出了一種融合 CNN 和 LSTM 的框架，將金融信息和股市歷史數據作為輸入，構建了七個不同的預測模型變種分類器。然而，使用深度學習預測股票價格或趨勢時，算法效果主要取決於預測準確度，並且在存在交易成本的情況下，高預測準確度並不完全代表最終收益率高，無法獲得由於股票交易活動而引起的未來懲罰或獎勵回報[4]。

區別於上述直接預測股票價格方法，基於強化學習方法將預測股價和投資動作結合在一起，直接以投資收益目標作為優化目標。當狀態與動作連續時，動作價值表空間過大，導致查表狀態和動作過於複雜，因而提出了使用神經網絡擬合狀態動作價值表，即深度強化學習。Deng 等[5]首次使用深度強化學習方法用於金融市場，深度學習自動感知動態市場條件，提取信息特徵，強化學習模塊與環境互動並做出交易決策，為進一步提高市場穩定性，引入模糊學習來減少輸入數據不確定性，實證結果較好。Li 等[6]將深度強化學習用於股票交易策略和股價預測，比較 DQN、Double DQN 和 Dueling DQN 三種不同的算法，均取得不錯的投資收益。Pendharkar 等[7]設計了基於在線策略和離線策略的離散狀態和動作模型以最大化投資回報和微分夏普比率。Jeong 等[8]針對交易數據不足的問題，提出遷移學習融合 Q-learning 處理高波動金融數據帶來的過擬合問題。Chakole 等[9]將強化學習的動作選擇與趨勢跟踪方法相結合，通過趨勢指標直接影響代理動作，結果顯示方法帶來較高的收益效果。但以上方法，在影響預期回報的情況下未降低策略風險，無法提供穩定的收益。

基於以上討論，本文提出了一種結合趨勢的深度強化學習股票交易模型尋找最佳交易策略。主要步驟如下：

- 1) 將股票開盤價和閉盤價相結合代表股票市場狀態。
- 2) 選取根據趨勢指標 RSI 指數調整後特定條件下的利潤作為獎勵函數。
- 3) 在股票市場中使用 DQN 做出交易決策。
- 4) 在股票數據集上的實驗結果顯示。

本文提出的模型評價指標均優於其他基準算法，使用趨勢指標調整獎勵函數明顯的有效改善模型的表現及效果。

## ● Problem statement :

股票交易是指在不斷變化的股票市場環境選擇股票並做出不同的交易動作從而改變資金在市場中的分配比例，最大化投資回報率並降低風險的過程。在強化學習交易模型中，通過深度強化學習網絡，以股票的歷史數據作為狀態，實現總收益最大化為目標，在每個交易時刻之前輸出一個交易動作，並為每個交易動作提出一個獎勵，通過自動交易將資金調整到最優，不斷進行計算和自我學習，從而實現優化的股票交易模型。

## ● Data description :

### 1. 實驗數據集

驗證本文提出的模型，在實驗中，從中證 100 指數中隨機選取 3 隻股票中信證券、保利發展和水泥進行實證分析，股票代碼如下：600030, 600048, 600585。交易數據是從 Tushare 金融社區下載的每日股價數據。實驗數據為訓練週期為 2010 年 1 月 1 日至 2017 年 12 月 31 日，測試期間從 2018 年 1 月 1 日至 2021 年 12 月 31 日。訓練期和測試期時長為 8a 和 4a。

### 2. 實驗環境

使用語言為 Python3.6.8，並採用 Pytorch1.10.1 為運行環境。

## ● Reinforcement learning and neural networks structure

狀態  $s$  由連續多個交易日開盤價和收盤價的變化情況構成。從市場中獲得每個交易日的股價波動信息，股價波動由兩部分組成，當前交易日開盤價較前一交易日閉盤價的變化率  $r_t^{oc}$  和當前交易日的閉盤價較當前交易日的開盤價的變化率  $r_t^{co}$ ，其中， $PO_t$  為  $t$  時刻股票的開盤價， $PC_t$  為  $t$  時刻股票的閉盤價。開盤價在一定程度上代表著市場消息面因素，收盤價則直接反映股價波動情況，這兩個變化率分

別代表著股票在休市和開市期間的股市信息。經實驗結果表明，當時間窗口 $T$ 取30時，模型效果最好。因此在本文中， $T = 30$ 。

$$s_t = (r_{t-T+1}^{oc}, r_{t-T+1}^{co}, r_{t-T+2}^{oc}, r_{t-T+2}^{co}, \dots, r_t^{oc}, r_t^{co}), T \geq 2$$

$$r_t^{oc} = \frac{PO_t - PC_{t-1}}{PC_{t-1}}$$

$$r_t^{co} = \frac{PC_t - PO_t}{PO_t}$$

趨勢使用技術分析指標來計算，比如相對強弱指數(Relative Strength Index, RSI)，順勢指標(Commodity Channel Index, CCI)等。如果趨勢向上，股價估計將上漲，相反，如果趨勢向下，股價估計將下跌，應結束之前的多頭頭寸。因中國股市不允許賣空，這裡不考慮空頭頭寸。在本文中，使用相對強弱係數RSI 決定趨勢，相對強弱係數是一個動量指標，它提供了一個超買或者超賣的信號。該指標的取值範圍為0 到100，如果其值低於30，表示超賣，其值高於70時，表示超買。相對強弱係數是根據股票前14個交易日的波動情況計算出來的，其中，average-gain 是14天內閉盤價上漲數之和的平均值，average-loss 是14天內閉盤價下跌數之和的平均值。為了規避風險，當市場處於超買情況時，股價極有可能發生反轉下跌，此時應抑制代理的買入和持有行為，鼓勵賣出行為，且下一天閉盤價下跌越多，對賣出行為的鼓勵越高；反之，當市場處於超賣情況時，股價極有可能發生反轉上漲，此時應抑制代理的賣出和持有行為，鼓勵買入行為，且下一天閉盤價上漲越多，對買入行為的鼓勵越高。通過調整特定條件下的獎勵函數值，造成對代理行為鼓勵或抑制的影響，這種調整通過在原先的獎勵值上乘以特定的影響係數 $m$ 來完成，特定條件如表1所示。

RSI<30			RSI>70		
動作 $a_t$	$r_{t+1}^c$	影響係數 $m$	動作 $a_t$	$r_{t+1}^c$	影響係數 $m$
1~10	0.05~0.1	1.3	1-10	-	0.8
1~10	0.03~0.05	1.2	0	-	0.9
1~10	0.01~0.03	1.1	-10~1	-0.01~0.03	1.1
0	-	0.9	-10~1	-0.03~0.05	1.2
-10~1	-	0.8	-10~1	-0.05~0.1	1.3

表1. 不同條件下調獎勵值的影響係數

$$RSI = 100 - \frac{100}{1 + \frac{\text{average} - \text{gain}}{\text{average} - \text{loss}}}$$

$$r_t = \begin{cases} m \times r_t^{profit}, & \text{滿足表一某條件} \\ m \times r_t^{profit}, & \text{其他} \end{cases}$$

## ● Simulations

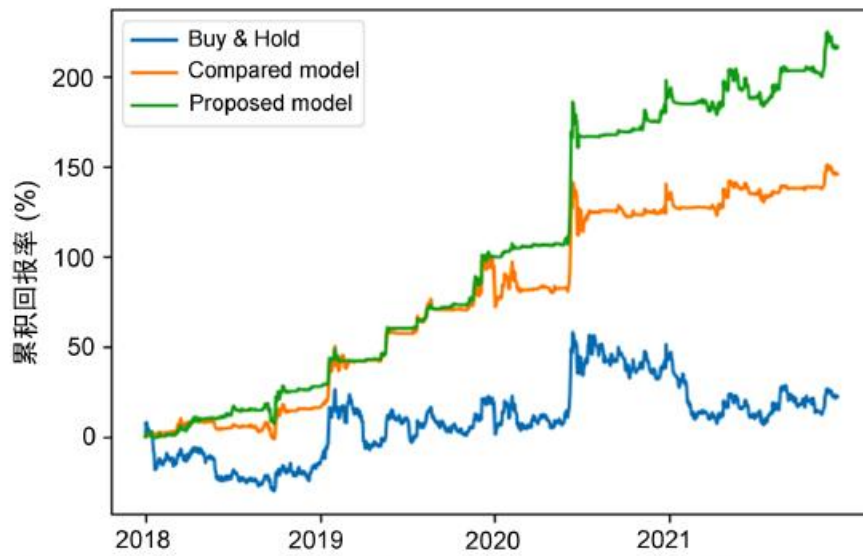


Figure 2. Cumulative returns of three investment models on 600030

圖 2. 中信證券上三種投資模型的累積收益率

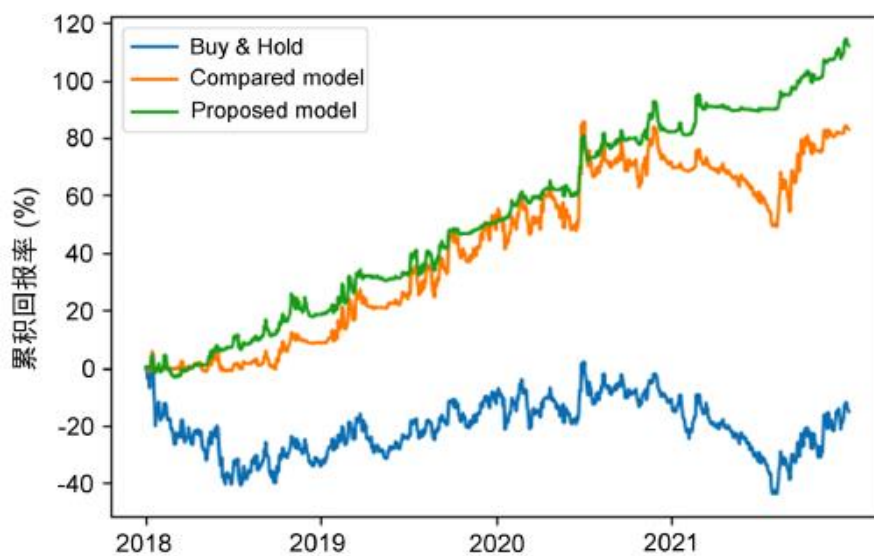


Figure 3. Cumulative returns of three investment models on 600048

圖 3. 保利發展上三種投資模型的累積收益率

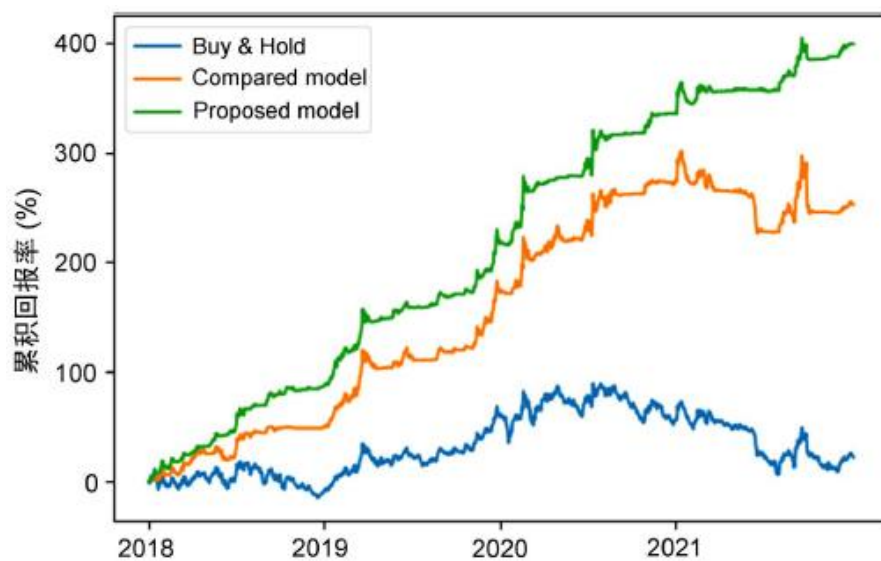


Figure 4. Cumulative returns of three investment models on 600585

圖 4. 海螺水泥上三種投資模型的累積收益率

投資模型	平均年回報(%)	年波動率(%)	夏普比率(%)
Buy & hold	5.1	35.32	19.62
Compared model	25.19	19.72	101.12
Proposed model	33.33	16.22	197.68

表 2. 中信證券上三種投資模型的評價指標對比

投資模型	平均年回報(%)	年波動率(%)	夏普比率(%)
Buy & hold	-4.04	37.46	-15.12
Compared model	16.31	19.8	121.89
Proposed model	20.65	12.9	514.62

表 3. 保利發展上三種投資模型的評價指標對比

投資模型	平均年回報(%)	年波動率(%)	夏普比率(%)
Buy & hold	-4.04	37.46	-15.12
Compared model	16.31	19.8	121.89
Proposed model	20.65	12.9	514.62

表 4. 海螺水泥上三種投資模型的評價指標對比

由上方表 2、表 3、表 4 展示了三種投資模型在中信證券、保利發展和海螺水泥上的評價指標對比情況。從表中可以看出基於傳統獎勵函數的強化學習投資模型可以獲得可觀的回報，並且其年波動率和夏普比率不錯。與前者相比，本文構建的投資模型的平均年回報率提升了 30%，年波動率降低 25%，夏普比率提升 50%以上。

## ● Conclusions

本文將股價趨勢與強化學習方法中的獎勵函數設定相結合，通過模型行動和股價在不同條件下的影響係數調整獎勵函數，使之構建成新的深度強化學習股票交易模型並應用於股票交易。本文在中國股票市場中選擇了 3 支股票進行投資實驗，實驗結果顯示，本文的模型表現優於其他對照組，在實驗期間的平均年回報更高，年波動率更低，且夏普比率更好，表明了在股票交易上的有效性，有較好的應用價值。但是本文的模型是基於一些假設進行的，不符合市場中實際投資者的投資方式。例如當交易量較大對股價造成影響時，本文的模型不適用，因此有待進一步研究與探索。

## ● References

- [1] Bao, W., Yue, J. and Rao, Y. (2017) A Deep Learning Framework for Financial Time Series Using Stacked Autoencoders and Long-Short Term Memory. PloS ONE,

12, e0180944.

- [2] Cai, S., Feng, X., Deng, Z., et al. (2018) Financial News Quantization and Stock Market Forecast Research Based on CNN and LSTM. International Conference on Smart Computing and Communication, Tokyo, 10-12 December 2018, 366-375.
- [3] Jiang, Z., Xu, D. and Liang, J. (2017) A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem. arXiv:1706.10059 [q-fin.CP].
- [4] Deng, Y., Bao, F., Kong, Y., et al. (2016) Deep Direct Reinforcement Learning for Financial Signal Representation and Trading. IEEE Transactions on neural Networks and Learning Systems, 28, 653-664.
- [5] Li, Y., Ni, P. and Chang, V. (2020) Application of Deep Reinforcement Learning in Stock Trading Strategies and Stock Forecasting. Computing, 102, 1305-1322.
- [6] Pendharkar, T. and Cusatis, P. (2018) Trading Financial Indices with Reinforcement Learning Agents. Expert Systems with Applications, 103, 1-13.
- [7] Jeong, G. and Kim, H.Y. (2019) Improving Financial Trading Decisions Using Deep Q-Learning: Predicting the Number of Shares, Action Strategies, and Transfer Learning. Expert Systems with Applications, 117, 125-138.
- [8] Chakole, J. and Kurhekar, M. (2020) Trend Following Deep Q-Learning Strategy for Stock Trading. Expert Systems, 37.
- [9] Leem, J. and Kim, H.Y. (2020) Action-Specialized Expert Ensemble Trading System with Extended Discrete Action Space Using Deep Reinforcement Learning. PLoS ONE, 15, e0236178.