**COMP5513**

**Financial Computing**

**Group Project Report**

**Group 15**

**21000018G Li Yat Long**

**21009566G Mak Ka Wai**

**21009856G Chan Siu Fung**

**21014947G Yu Hin Chung Nikko**

**Table of Contents**

## 1. Paper Studies

## 1.1 Li Yat Long's Paper – Stock Market's Price Movement Prediction with LSTM Neural Networks [1]



Fig. 1 Screenshot of the paper



Fig. 2 Methodology of the paper [1]

Studies on stock prices predictions methods have been investigated for decades due to its complexity and dynamism, where Machine Learning is one of the common approaches among them. The recurrent Long-Short term memory (LSTM) network which proves effective on Natural

Language Processing (NLP), was introduced to stock prices predictions in this paper. Compared to traditional feed-forwarded networks, LSTM allows neurons passing in multiple directions with weights given depending on the age of the data, forgetting memory when the data is evaluated as irrelevant and better performance for long data sequence input.

The historic price data from 2008 to 2015 were collected from the BM&F Bovespa stock exchange with the format of time series of candles including open, close, high, low and volume in a granularity of 15 minutes. After the data were collected, both log-return transformation and exponential smoothing were performed in order to stabilize the mean and variance of the time series, as well as reduce the random and noise on the price series. The data would be further transferred to the technical indicator using the TA-Lib library, generating 175 values which represent the future price, trading volume, movement tendency, visual graphical patterns of a stock in each period, by determining or predicting the characteristics of data. Finally, tuples of price data and technical indicators would be applied to train the neural network to predict if the closing price of the next period will be higher than the current closing price. The neural network would be trained at the end of each trading day to prepare for the prediction on the next day, by using the latest sets of data from the last 10 months till the current day with updated weights.

The LSTM model was compared to another three traditional methods including multi-Layer perception, Random Forest, and pseudo-random model with the same data input. The results showed that LSTM performed much better in terms of average accuracy, with 1 to 6 % better results. In addition, LSTM have achieved a high return ratio with significantly less risk compared to other strategies. However, the standard deviation was obviously larger in most cases compared to multi-Layer perception and Random Forest, which was the area the paper suggested to be improved.

## 1.2 Mak Ka Wai's Paper – Stock market analysis using candlestick regression and market trend prediction (CKRM) [2]



Fig. 3 Screenshot of the paper

In this paper, it proposed a stock prediction system which hast three main modules. The first module is trendline identification consisting of MACD and candlestick graph to indicate buy and sell signal. The second module is using K Nearest Neighbor algorithm to predict the stock prices. Features include MACD, RSI, Bollinger bands and candle stick pattern detection. The last module is pattern formation such as formed by Open, High, Low, Close (OHLC) and RSI over past few days to predict the trend of next day. Also, the proposed system is implemented Django web framework and dataset is using Flask REST API which is getting live data in an open-source platform called Quandl.

Fig. 4 Proposed system of the paper [2]

Moreover, this paper has compared the prediction accuracy from three common algorithms which are KNN regression, Support Vector Machine (SVM) and Linear regression to predict the stock prices of different companies (TCS, HCL, Oracle Finance, MindTree, NIITech, WiPro). Features include MACD, RSI, Bollinger bands, candle stick pattern detection. Moreover, this paper has used exponential moving average (EMA) and relative strength index for features and variation to have higher prediction for stocks.

| Company/machine learning algorithms | kNN regression | SVM | Linear regression |
|---|---|---|---|
| TCS | 2948.679 | 2930.067 | 1962.287 |
| Infosys | 1149.887 | 1181.541 | 1921.568 |
| Wipro | 302.4868 | 443.0676 | 481.3753 |
| HCL | 930.1876 | 839.772 | 839.011 |
| Oracle Finance | 3882.339 | 4046.162 | 3192.046 |
| MindTree | 785.0703 | 730.844 | 686.5661 |
| Hexaware | 353.4225 | 319.8947 | 323.8736 |
| Tech Mahendra | 613.0417 | 633.1691 | 664.855 |
| NIITech | 850.5223 | 704.8029 | 398.9765 |

Fig. 5 Prediction value of selected companies [2]

Below is the prediction accuracy from the paper. It points out that KNN regression has the best performance. KNN has 90%- 96% accuracy compared to Linear regression (80% - 95%) and Support vector machine (60% - 80%).

**Table 2** Accuracy range of regression algorithms

| Method | Accuracy range (in %age) |
| --- | --- |
| k-NN regression | 90–96 |
| Linear regression | 80–95 |
| Support vector machine | 60–80 |

Fig. 6 Finding of the paper

## 1.3 Chan Siu Fung's Paper – Two-channel Attention Mechanism Fusion Model of Stock Price Prediction Based on CNN-LSTM [3]

**Two-channel Attention Mechanism Fusion Model of Stock Price Prediction Based on CNN-LSTM**

LIN SUN, WENZHENG XU, and JIMIN LIU, College of Intelligent Equipment, Shandong University of Science and Technology, Tai'an, China

Using hierarchical CNN, the company's multiple news is characterized as three levels: sentence vectors, chapter vectors, and enterprise sentiment vectors. By combining the stock price data with the news lyric data at the same time, the influence of news on price is used to achieve correlation analysis of news information and stock prices. A two-channel attention mechanism fusion model based on CNN-LSTM is proposed. After the dual-channel feature extraction, the attention layer fusion layer is used to convert the weighted values of LSTM hidden variables, so the stock price can be predicted with the news text.

CCS Concepts: • **Theory of computation** → *Automated reasoning;* • **Networks** → *Network performance analysis;*

Additional Key Words and Phrases: CNN-LSTM, stock prediction, attention mechanism, two-channel

**ACM Reference format:**
Lin Sun, Wenzheng Xu, and Jimin Liu. 2021. Two-channel Attention Mechanism Fusion Model of Stock Price Prediction Based on CNN-LSTM. ACM Trans. Asian Low-Resour. Lang. Inf. Process. 20, 5, Article 83 (July 2021), 12 pages.
http://dx.doi.org/10.1145/3453693

83

## 1 INTRODUCTION

With the development of machine learning and deep learning algorithms, many scholars use RNN, CNN, LSTM, seq2seq, and attention mechanism to predict stock prices. In related researches, RNN has advantages in processing time series due to its memory ability before and after the sequence. It is used by the most researchers. RNN-based LSTM, seq2seq, and attention mechanism models are used as well. Sun Xiang proposed a secondary neural network structure that combines **recurrent neural network (RNN)** with kernel feature extraction for stock price forecasting.[1] Ren Jun [2] proposed the regularized long-term and short-term memory neural network LSTM model applied to Dow Jones index prediction. Cho [3] and Sutskever [4] have successively proposed the SequencetoSequence model, which implements time series variable length input and output prediction through decoder and encoder, also known as "Encoder-Decoder" model.

Authors' addresses: X. Wenzheng, S. Lin, and L. Jimin, No. 233 Daizong Street, Tai'an, Shandong Province, 271019, China; emails: wenzhengxu@yandex.com, Lin Sun 128735221@qq.com, Jimin Liu j.liu0291@163.com.
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2021 Association for Computing Machinery.
2375-4699/2021/07-ART83 $15.00
http://dx.doi.org/10.1145/3453693

ACM Trans. Asian Low-Resour. Lang. Inf. Process., Vol. 20, No. 5, Article 83. Publication date: July 2021.

Fig. 7 Screenshot of the paper

Fig. 8 Two-channel attention mechanism fusion CNN-LSTM Model [3]

Characteristics of proposed model

A dual channel attention mechanism fusion model combining both stock news and stock prices was introduced in this paper for stock price prediction.

In the model, the stock data and news corpus are processed separately in two different channels before the attention fusion layer that can avoid mutual interference between the stock price features and stock news features. A bidirectional LSTM model is used to process the stock prices while a CNN model is used to extract data from stock news.

To capture the semantic information related to the rise and fall of stocks from the long news items more accurately, the news is combined and analyzed together. News is segmented and words are stopped and embedded to vectors in embedding layer, the vectors are further compressed in three layers into different levels namely sentence, chapter, enterprise levels in the CNN model.

The fusion layer combines the features from both channels and then passes the output to MLP model to obtain final output. Through the fusion layer, a link between the corporate news and the stock price value can be established.

The feature and information fusion focusing on the derived attributes and data enhances the intrinsic perspective in stock market prediction. Meanwhile, the model fusion combines both the ability of CNN to express text space features and the ability of Bi-LSTM to express time dimension through attention mechanism, so this dual-channel fusion deep learning model has the ability to express time and text space of stock price and stock news data simultaneously.

Comparison of models

The training set loss value and test set loss value of the model were compared with those of purely Bi-LSTM model and single-channel CNN-LSTM in terms of the average value of MSE, RSME and MAE of each stock. The stock data was derived from the closing price of 411 stocks in the period of January 2016 to December 2017 obtained by the tushare interface. The news data were extracted from over 23000 news items related to those stocks within that period from dozens of financial websites like Securities Times, China Securities Network, Sina Finance, and China Economic Net. The result showed that the proposed model has the least loss among three. This implies that the proposed model can improve the accuracy of stock price prediction.

Return from trading based on stock price prediction by the model

20 stocks were randomly selected and being forecasted with the model, trading was done by throwing at peak, purchasing at bottom, the average income of selected stocks was found. The total stock return is 290%, the annual compound yield is as high as 126.5%, the positive income ratio is 64% which far exceeds the benchmark CSI 300 index yield. This showed that the predictions based on this model can achieve sustained positive returns with good predictive results.

## 1.4 Yu Hin Chung Nikko's Paper – Stock Closing Price Prediction Using Machine Learning Techniques [4]



Fig. 9 Screenshot of the paper

This paper adopts two machine learning methods to predict the stock price, namely Random Forest (RF) and artificial neural networks (ANN). There are several parameters they had made to make predictions, including:

- Stock High Price – Stock Low Price (H-L)
- Stock Close Price – Open Price (O-C)
- 7,14,21 Moving Average
- 7 days standard deviation

This paper uses stocks from companies like Nike, Goldman Sachs and JP Morgan for training and testing. As for evaluating the model performance, the paper uses Root Mean Square Error (RMSE), Mean Absolute Percentage Error (MAPE) and Mean Bias Error (MBE).

| Company | ANN | | | RF | | |
|---|---|---|---|---|---|---|
| | RMSE | MAPE | MBE | RMSE | MAPE | MBE |
| Nike | 1.10 | 1.07% | -0.0522 | 1.29 | 1.14% | -0.0521 |
| Goldman Sachs | 3.30 | 1.09% | 0.0762 | 3.40 | 1.01% | 0.0761 |
| JP Morgan and Co. | 1.28 | 0.89% | -0.0310 | 1.41 | 0.93% | -0.0313 |
| Johnson & Johnson | 1.54 | 0.70% | -0.0138 | 1.53 | 0.75% | -0.0138 |
| Pfizer Inc. | 0.42 | 0.77% | -0.0156 | 0.43 | 0.8% | -0.0155 |

Fig. 10 Paper results. [4]

The results shown that ANN generally behaves better than RF. For instance, using the Pfizer dataset, ANN has a score of 0.42 (RMSE) ,0.77 (MAPE) and –0.0156 (MBE), while RF has 0.43 (RMSE), 0.8 (MAPE) and –0.0155 (MBE).

This paper is part of our inspiration to develop this model as the orange software include ANN and RF models, as well as some of the parameters such as moving averages and the standard deviation.

## 2. Stock price prediction

### 2.1 Feature selection

The 2800 stock price data from 2011 to 2021 were downloaded from the internet and features were selected and calculated from the raw data, these include:

1. Close price

2. Close price 1 to 5 days before the current date

3. Percentage change of price 1 to 5 days ago compared to today's price

4. 10-, 20-, 100- & 200-day moving average (MA) which calculates the average price of recent 10/ 20/ 100/ 200 days

5. Uptrend indicated when the closing price is higher than the 10-, 20-, 100- & 200-day moving averages or downtrend indicated when the closing price is lower than the 10-, 20-, 100- & 200-day moving averages

6. Buy signal indicated when the closing price moves over 10-, 20-, 100- & 200-day moving averages or sell signal when the closing price moves below 10-, 20-, 100- & 200-day moving averages

7. Buy signal indicated when the closing price above the moving averages rebounds upwards after touching the moving averages or sell signal indicated when the closing price below the moving averages rebounds downwards after touching the moving averages

8. Gain/ Loss in closing price compared to last trading day

9. Average Gain/ Loss of 14 days

10. Relative Strength – Average gain divided by average loss

11. Relative Strength Index

12. Buy signal indicated when the RSI leaves the oversold region or sell signal indicated when the RSI leaves the overbought region

13. 12- & 26-day exponential moving average (EMA)

14. MACD – Difference between 12-day EMA and 26-day EMA

15. Signal – 9-day EMA of MACD line

16. Histogram – MACD value minus Signal value

17. Simplified signal line crossover – Buy signal indicated when MACD line just goes over the signal line or sell signal indicated when MACD line just goes below the signal line

18. Lower/ Upper Bollinger band

19. Simplified mean reversion – Buy signal indicated when closing price below the lower Bollinger Band goes over it or sell signal indicated when closing price over the upper Bollinger Band goes below it



Fig. 11 Excel files for calculating the features
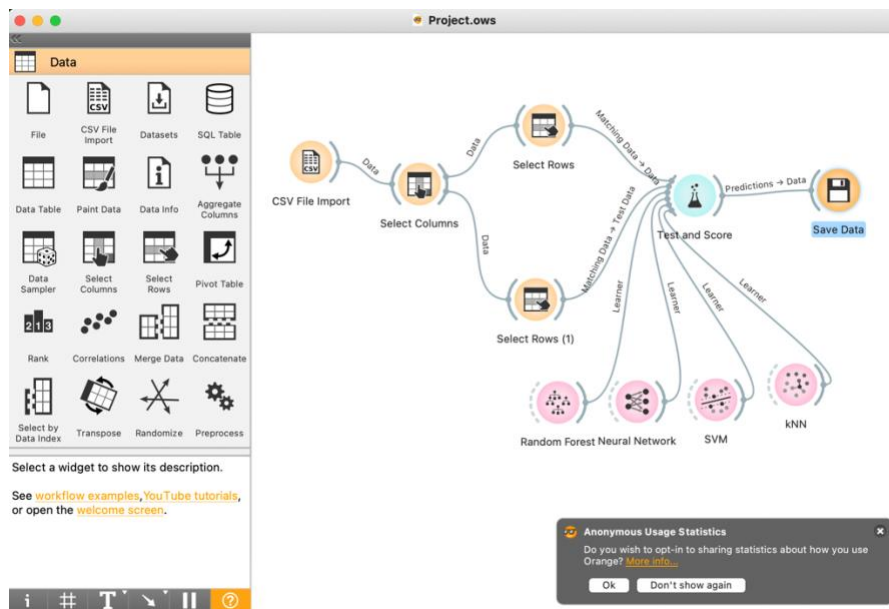
## 2.2 Machine Learning Model Training



Fig. 12 Orange training Interface

The software *Orange* was selected to train and build the machine learning models for predicting the stock price, which includes:

- Random Forest – it predicts the results by comparing the input with the root of the decision tree. The predicted price is the mean of the past prices which lie in the region that are larger / smaller than the root node of the branch.

- Neural Network – past prices are applied through interconnected hidden layer network with nodes activated by the value of threshold which is calculated by the input of features and the weight to train a model. By inputting the new values to the trained neural network in input layer, the values will be classified in output layer through activating the nodes in hidden layer.

- Support Vector Machine (SVM) – price is classified by finding a line separating the data into two classes which the margin by comparing with the support vector in both classes should be the largest.

- k Nearest Neighbor (kNN) – features are inputted, the square root of the sum of the squares of the difference between the feature of a test datapoint and that of a point in the model is

calculated as the distance and the mean or weighted mean of the top k-closest datapoint will be the predicted value.

On the other hand, the features selected on section 2.1 were spilt into 5 sets for training, producing 5 sets of stock price prediction model, which includes:

A. Close Price, Close price 1 to 5 days before
B. Close Price, Percentage change of price 1 to 5 days ago compared to today's price
C. Close Price, 10 day moving average, RSI, MACD, Signal, Histogram
D. Close Price, RSI
E. Close Price, MACD, Signal, Histogram

The above feature sets A-E were applied to 3 test cases in Orange, generating 60 price prediction model for comparison.

- Case 1: use price from year 2011 – 2020 as training data
- Case 2: use price from year 2014 – 2020 as training data
- Case 3: use price from year 2016 – 2020 as training data

### 2.3 Training Results

Below is the training result of case 1 which used Year 2011 – 2020 historical data to predict the stock price of 2021.

| Case 1 (Year 2011 – 2020 as training data) | | | | |
|---|---|---|---|---|
| Feature set | MAE (Random Forest) | MAE (Neural Network) | MAE (SVM) | MAE (KNN) |
| A | 0.304 | 0.286 | 1.066 | 0.297 |
| B | 0.771 | 0.759 | 0.800 | 0.751 |
| C | 0.339 | 0.367 | 0.345 | 0.347 |
| D | 0.302 | 0.309 | 0.555 | 0.333 |
| E | 0.339 | 0.410 | 0.426 | 0.295 |

Fig. 13 Training results of Case 1

Below is the bar chart of training results of Case 1. We found that feature A has the lowest MAE. For instance, Neural Networks of Feature A got 0.286 MAE which is the lowest result of training results of Case 1.
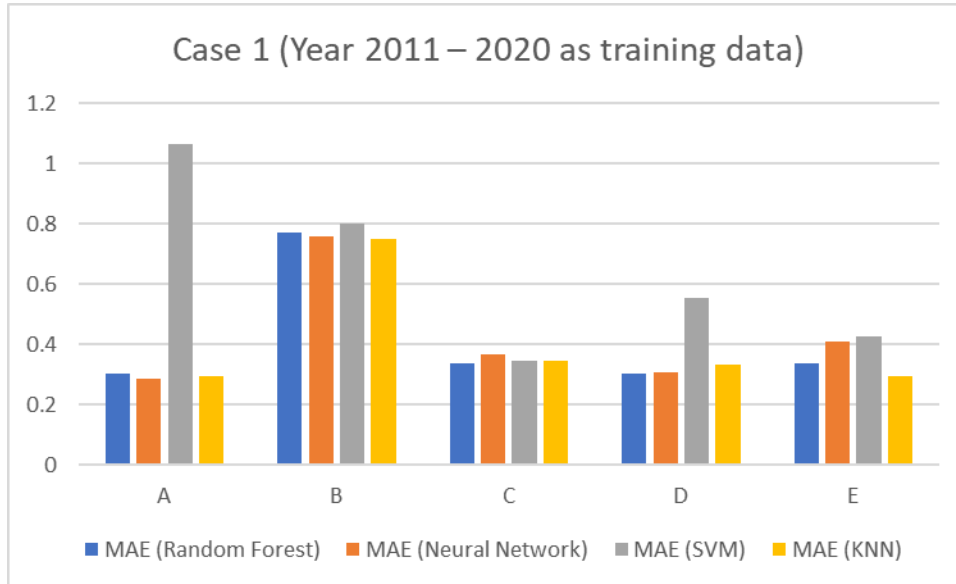


Fig. 14 Training results of Case 1 (bar chart)

Below is the training result of case 2 which used 2014 – 2020 historical data for training.

| Feature set | MAE (Random Forest) | MAE (Neural Network) | MAE (SVM) | MAE (KNN) |
|---|---|---|---|---|
| A | 0.306 | 0.351 | 0.570 | 0.301 |
| B | 0.271 | 0.271 | 0.431 | 0.280 |
| C | 0.328 | 0.513 | 0.341 | 0.346 |
| D | 0.304 | 0.360 | 0.543 | 0.333 |
| E | 0.341 | 0.542 | 0.562 | 0.305 |

Case 2 (Year 2014 – 2020 as training data)

Fig. 15 Training results of Case 2

Below is the bar chart of training results of Case 2. We found that feature B has the lowest MAE. For instance, Neural Networks and Random Forest of Feature B got 0.271 MAE which is the lowest result of training results of Case 2. It showed that Case 1 and Case 2 have different outcomes in different features.
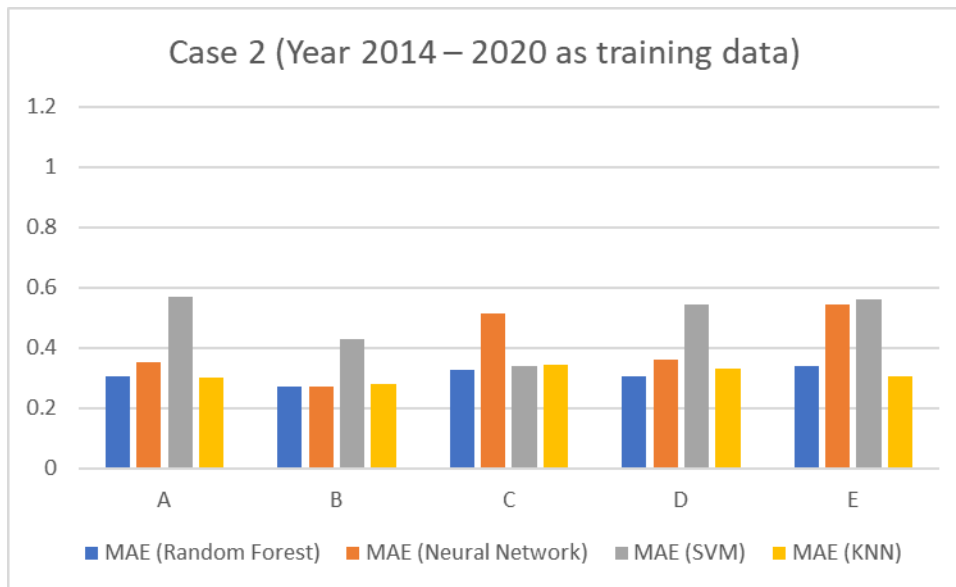


Fig. 16 Training results of Case 2 (bar chart)

Below is the training result of case 3 which used historical data within 2016 – 2020 for training our model.

Case 3 (Year 2016 – 2020 as training data)

| Feature set | MAE (Random Forest) | MAE (Neural Network) | MAE (SVM) | MAE (KNN) |
|---|---|---|---|---|
| A | 0.310 | 0.471 | 1.016 | 0.305 |
| B | 0.285 | 0.272 | 0.329 | 0.279 |
| C | 0.338 | 0.997 | 0.328 | 0.366 |
| D | 0.301 | 0.525 | 0.763 | 0.364 |
| E | 0.342 | 0.912 | 0.527 | 0.307 |

Fig. 17 Training results of Case 3

Below is the bar chart of training results of Case 3. We found that feature B has the lowest MAE. For instance, Neural Networks of Feature B got 0.272 MAE which is the lowest result of training results of Case 3. It showed that Case 2 and Case 3 have same training result for predicting 2800 stock price in 2021.



Fig. 18 Training results of Case 3 (bar chart)

In sum up of training phase, we found out that feature B (Close Price, Percentage change of price 1 to 5 days ago compared to today's price) has highest accuracy in Case 2 and Case 3. For instance, Neural Network algorithm has 0.271 MAE in Case 2 and 0.272 MAE in Case 3. We believe Feature B is a better approach for predicting stock based on our studies.

## 3. Trading Strategies

The predicted stock price obtained from Case 3 feature set B were used for evaluating the effectiveness of the trading strategies. Total 9 strategies were developed (strategy 1A, 2A, 2B, 3A, 3B, 4A, 4B, 5A & 5B) and $1000000 initial cash were applied to all strategies to simulate stock bought/ sold in year 2016 and 2011.

An assumption had been made to all strategies with A suffix, which the fixed daily transaction amount approach was applied:

- Max transaction amount = Total cash/Number of trade days in the year
- Max number of shares to buy (Nmax) = Max transaction amount/Close price of the day

On the other hand, for all strategies with B suffix, fixed percentage approach was applied which bought with fixed percentage of cash and sold fixed percentage of shares.

### 3.1 Trading Strategy 1 (1A)

For strategy 1A, Nmax of shares are purchased each day regardless the stock is rise or fall.

### 3.2 Trading Strategy 2 – Purely based on price prediction

For strategy 2A, Nmax of shares are purchased when predicted price rise and Nmax shares are sold when predicted price fall. No action when there no price change in prediction.

For strategy 2B, 50% of the remaining capital are used to purchase shares when rise is predicted and 50% of the shares are sold when fall is predicted. No action when there no price change in prediction.

### 3.3 Trading Strategy 3 – Partially based on price prediction, partially based on technical indicator signals

For strategy 3A, half of Nmax shares is purchased when rise predicted and half of Nmax shares is sold when fall predicted. In addition, another half of Nmax shares is purchased when there is buy signal which can be indicated by moving average cross-over, RSI overbought/oversold, MACD cross-over or mean reversion, and half of Nmax shares is sold when there is sell signal.

For strategy 3B, 25% of capital is used to purchased stocks when rise predicted and 25% of stocks is sold when fall predicted. In addition, another 25% of capital is used to purchased stocks when there is buy signal and 25% of stocks is sold when there is sell signal.

### 3.4 Trading Strategy 4 – Purely based on technical indicator signals

For strategy 4A, Nmax of shares is purchased when there are buy signal and Nmax of shares is sold when there are sell signal.

For strategy 4B, 50% of capital is used to purchase shares when there are buy signal and 50% of shares is sold when there are sell signal.

### 3.5 Trading Strategy 5 – Purely based on 4 moving averages (10-Day MA, 20-Day MA, 100-Day MA, 200-Day MA)

For strategy 5A, quarter of Nmax shares are purchased when the close price is higher than the 10-day MA and quarter of Nmax shares are sold when the close price is lower than the 10-day MA.

The same practice also applied to the remained 3 quarters of Nmax shares with 20-Day MA, 100-Day MA, 200-Day MA and the buy or sell signal.

For strategy 5B, 1/10 of capital are used to purchase shares when the close price is higher than the 10-day MA and 1/10 of shares ar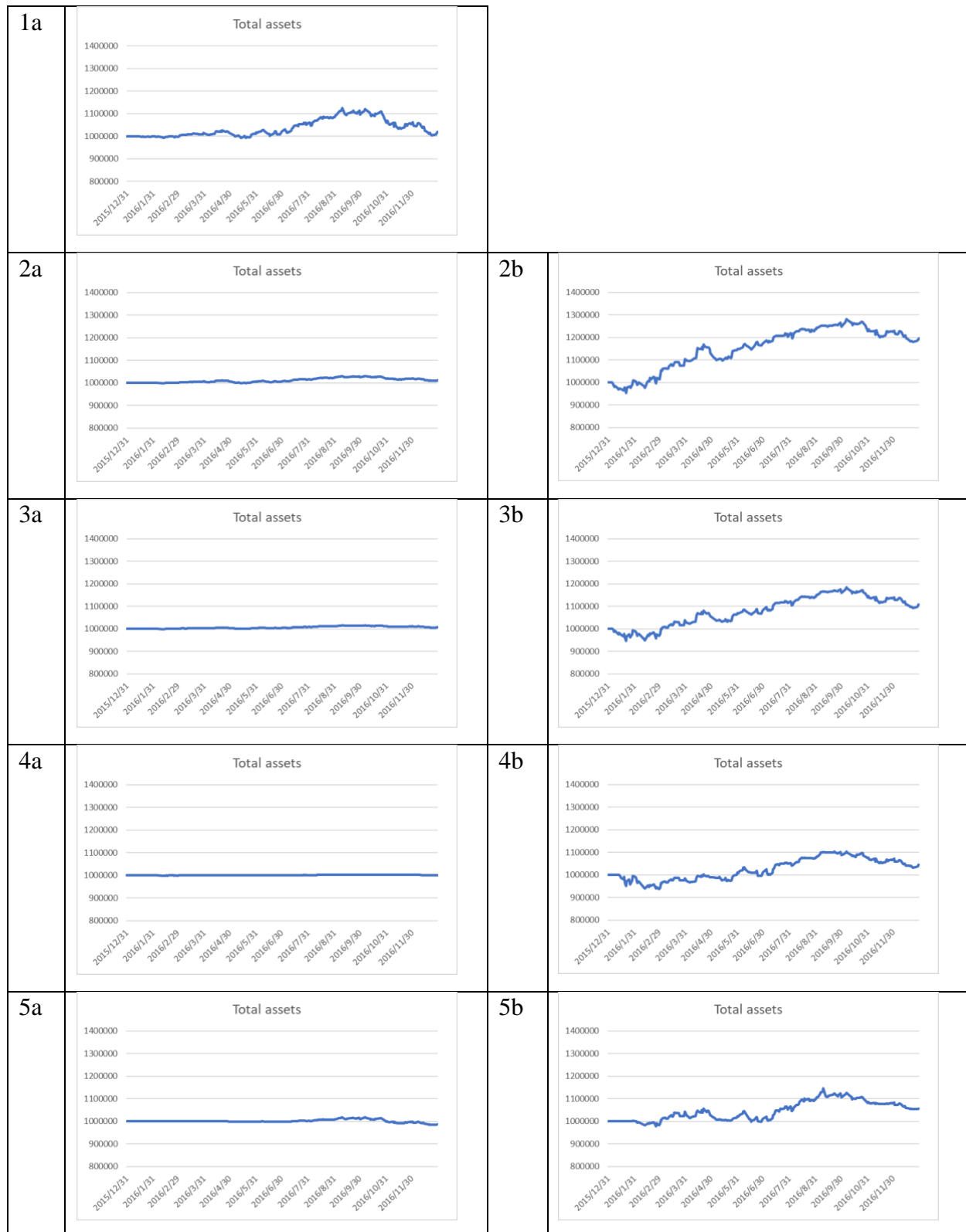e sold when the close price is lower than the 10-day MA. In addition, 1/20 of capital are used to purchase shares when the close price is higher than the 20-day MA and 1/20 of shares are sold when the close price is lower than the 20-day MA. Furthermore, 1/100 of capital are used to purchase shares when the close price is higher than the 100-day MA and 1/100 of shares are sold when the close price is lower than the 100-day MA. Finally, 1/200 of capital are used to purchase shares when the close price is higher than the 200-day MA and 1/200 of shares are sold when the close price is lower than the 200-day MA.

## 3.6 Results & Comparison



Fig. 19 Variation of total assets (2016)

Fig. 20 Variation of total assets (2021)

| Trading strategy | 2016 | 2021 |
|---|---|---|
| 1 | 1018685 (1.87%) | 863713.8 (-13.63%) |
| 2A | 1011811 (1.18%) | 994956.6 (-0.5%) |
| 2B | 1196393 (19.64%) | 898953.9 (-10.1%) |
| 3A | 1007183 (0.72%) | 996909.0 (-0.31%) |
| 3B | 1107773 (10.78%) | 936439.3 (-6.36%) |
| 4A | 1001364 (0.14%) | 997823.6 (-0.22%) |
| 4B | 1045212 (4.52%) | 953520.3 (-4.65%) |
| 5A | 987164.2 (-1.28%) | 975545.8 (-2.45%) |
| 5B | 1055626 (5.56%) | 884471.5 (-11.55%) |

Fig. 21 Results of all trading strategy

Among all trading strategies, strategy 2B which based on purely price prediction, achieved better profits than other strategies in bull market – 2016, with 19.64% increase in capital. Strategy 3B also did well in 2016 with 10.78% gains in capital. In contrast, strategy 5A did the worst with -1.28% loss in capital.

However, when it was bear market (year 2021), strategy 4B achieved the least loss with only 0.22% loss in capital. Strategy 3B, and 2B also did well with 0.31% and 0.5% loss. In contrast, strategy 5B and 1 did the worst with 11.55% and 13.63% loss.

For buyers who consider profit is the most important and have the ability accept larger loss, strategy 2B are suggested which performs the best in bull market while having larger loss than strategy 3B/4B in bear market. On the other hand, for the buyers who consider defence the most, strategy 4B is suggested which performs the best in bear market while having 0.14% gain in capital in bull market.

Meanwhile, strategy 1 and 5 are not suggested for any buyers:

- For strategy 1, it achieved the worst in the bull market and only gain 1.87% increase in bull market.

- For strategy 5, 5A is the only strategy obtain loss in bull market and 5B is obtain a huge loss in bear market.

## References

[1]    D. M. Q. Nelson, A. C. M. Pereira, and R. A. d. Oliveira, "Stock Market's Price Movement Prediction With LSTM Neural Networks," presented at 2017 Stock Market's Price Movement Prediction With LSTM Neural Networks. [Online]. doi: 10.1109/IJCNN.2017.7966019 [Accessed Feb 14, 2022].

[2]    M. Ananthi, and K. Vijayakumar, "Stock market analysis using candlestick regression and market trend prediction (CKRM)," *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, pp. 4819-4826, 2021. [Online]. doi: 10.1007/s12652-020-01892-5 [Accessed Feb 16, 2022].

[3]    L. Sun, W. Xu, and J. Liu, "Two-channel Attention Mechanism Fusion Model of Stock Price Prediction Based on CNN-LSTM," *ACM Trans. Asian Low-Resour. Lang. Inf. Process. 20*, vol. 5, 2021. [Online]. doi: 10.1145/3453693 [Accessed Feb 18, 2022].

[4]    M. Vijh, D. Chandola, V. A. Tikkiwal, and A. Kumar, "Stock Closing Price Prediction Using Machine Learning Techniques," *Procedia Computer Science*, vol. 167, pp. 599-606, 2020. [Online]. doi: 10.1016/j.procs.2020.03.326 [Accessed Feb 18, 2022].

**Appendix**

    **Individual Report – 21000018G Li Yat Long**

In this project, I oversee feature selection, model training, calculating MAE on case 3 feature set B, verifying results of trading strategies, and presentation slides & report writing.

During feature selection, different combinations of feature sets were discussed, and feature set A and B were selected to follow the method in the lab to ensure its performance while offering an opportunity to compare the results with other feature sets. Feature sets C tried included all indicators learned during lecture and feature sets D and E only included one type of indicators to compare the performance with feature set C and the methods learned from lab (feature sets A and B).



Fig. 22 Model training by Orange

In the training process, I have assisted on inserting the CSV file to Orange, selected the appropriate training and testing data, and obtained the results from Test and Score module. All feature sets except feature B can be directly concluded from Orange. The MAE results obtained from Orange in feature set B calculated by the percentage change in price, which was different from the MAE

of actual price our team would like to obtain. Therefore, I assisted in calculating the MAE in feature sets B following the methods learned from lab.

On the other hand, while implementing trading strategies, I assisted in verifying the excel table made by my groupmates. This was an important step not only to verify the correctness of the results but also to improve my understanding of the implementation process and the difference between the variety of strategies made by my groupmates.



Fig. 23 Calculate the actual MAE of feature sets B

Fig. 24 Verifying the results of trading strategies

Finally, during the report and presentation slides writing, I discussed with my teammates to include which kind of graphs or table to represent our results clearly to our audience. One of the results was the chart shown in Fig. 23. In addition, I have also implemented both section 2 and 3 of report and presentation slides partly. Through this process, I have understood more of the parts not implemented by me like trading strategies and have better understanding of the results on both section 2 and 3. For section 1, I have understood more on the stock price prediction using LSTM which is not mentioned in the lectures, with is methodology and mechanism.

**Individual Report – 21009566G Mak Ka Wai**

In this project, I am responsible for feature selection discussion, prediction tools discussion and trading strategy discussion, model training, calculating MAE of feature set B in case 1 and case 2, writing presentation slides & report and presentation.

Firstly, I have discussed with team members how to predict 2800 stock price. For instance, what tools we are planning to use. Which years are we predicting? Which features will we use to train our algorithm? After discussion, we decided to use Orange to train our dataset. Also, we will have 5 set features shown in reports to see which feature will get the highest accuracy. Moreover, we set 3 different cases to predict 2800 stock price in 2021.

Apart from discussion, I have helped to use Orange to train the model for case 1 and case 2 to find out the finding of our research. After finding out the result of each feature, I have exported the finding as excel and marked down in our report and presentation slide. Also, I have calculated MAE of feature B in excel based on the instruction of the workshop.



| | | | |
|---|---|---|---|
| | 2800.HK(RemovedNA)_ModelB.csv | March 19 | YU, Nikko [Student] |
| | 2800.HK(RemovedNA)_ModelB.xlsx | March 27 | MAK, kawai [Stude… |
| ○ | 2800.HK(RemovedNA)_ModelC.csv  … | March 19 | YU, Nikko [Student] |
| | 2800.HK(RemovedNA)_ModelD.csv | March 19 | YU, Nikko [Student] |
| | Case 1 Feature A.xlsx | March 27 | CHAN, Current [Stu… |
| | Case 1 Feature B.xlsx | March 24 | MAK, kawai [Stude… |
| | Case 1 Feature C.xlsx | March 28 | YU, Nikko [Student] |
| | Case 1 Feature D.xlsx | March 23 | MAK, kawai [Stude… |
| | Case 1 Feature E.xlsx | March 23 | MAK, kawai [Stude… |
| | Case 2 Feature A.xlsx | March 23 | MAK, kawai [Stude… |
| | Case 2 Feature B.xlsx | March 24 | MAK, kawai [Stude… |
| | Case 2 Feature C.xlsx | March 23 | MAK, kawai [Stude… |
| | Case 2 Feature D.xlsx | March 23 | MAK, kawai [Stude… |
| | Case 2 Feature E.xlsx | March 23 | MAK, kawai [Stude… |
| | Feature E (1).xlsx | March 23 | CHAN, Current [Stu… |
| | Project_feature_a.xlsx | March 19 | YU, Nikko [Student] |

Fig. 25 Training result of case 1 and case 2
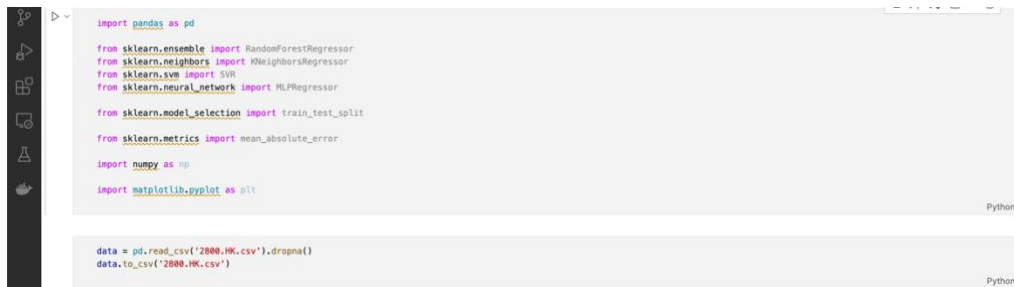
Fig. 26 Calculating MAE of feature B

Apart from training models, I have discussed trading strategy and graphs of our findings in the presentation slide with my team members. For instance, I am responsible for presenting the training results during the presentation.

After working on this project, I have learnt how to analysis and predict a stock using machine learning algorithms in real world stock market. This project enhanced my confidence in investing in stocks or building up my financial portfolio. Moreover, I think Orange is easy to implement different machine learning algorithms to predict stocks in the market.

**Individual Report – 21009856G Chan Siu Fung**

<u>Contribution</u>

During the project, I have mainly contributed for the calculation-required tasks and graph-plotting. I have removed the rows with empty values from the raw historical data that my groupmates have downloaded from the Internet by using the skill learnt in the extra Python tutorial.

```python
import pandas as pd

from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import SVR
from sklearn.neural_network import MLPRegressor

from sklearn.model_selection import train_test_split

from sklearn.metrics import mean_absolute_error

import numpy as np

import matplotlib.pyplot as plt
```

```python
data = pd.read_csv('2800.HK.csv').dropna()
data.to_csv('2800.HK.csv')
```

Fig. 27 Python codes for removing empty entries in the csv files downloaded from the Internet

I have calculated the technical indicators and prepared the derived features by using Excel.

Fig. 28 Screenshot of derived features of technical indicators

I have calculated the next price change percentage, the moving averages (10-day, 20-day, 100-day, 200-day), found if it is an upward trend or downward trend by comparing the moving averages with the close price, the buy/sell signals detected by looking for the cross-over between close price line and moving average lines or the pull-back from the moving average lines.

I have also calculated the daily gain, daily loss, average gain, average loss, relative strength, and relative strength index, looked for buy/sell signals due to bullish or bearish divergence.

I have then calculated the 12-day EMA and 26-day EMA for calculating the MACD and signal, and thus the histogram and the cross-over between MACD line and signal line.

I have also calculated the lower and upper Bollinger bands and looked for mean reversion. In the trading strategies part, I have proposed different strategies, actualized the trading and calculated the final total assets for each strategy.

I have plotted clustered column bar charts for illustrating the training error of different models trained with different sets of features in different cases and the plots for the daily variation of total assets with different trading strategies.



Fig. 29 Clustered column bar chart for comparing training errors of different models trained with different sets of features in Case 1



Fig. 30 Daily variation of total assets with different trading strategies

Learning reflection
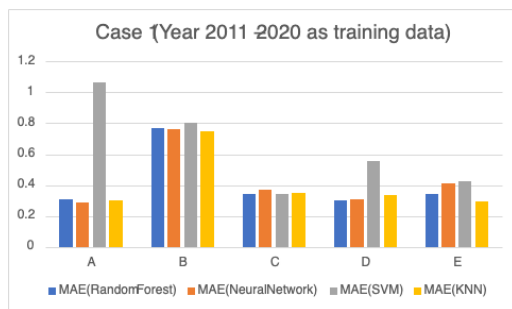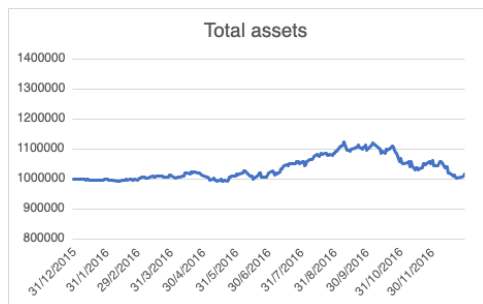
Through the study of paper, I think I can understand the paper with what I have learnt in the lectures in both this course and the course of Machine Learning and Data Analytics, as there are so many

technical terms in the Machine Learning aspect (terms like features, MAE, MSE, RMSE, MLP, CNN, RNN, LSTM, sigmoid function, tanh function, hidden states, softmax function, etc.) and the Finance aspect (terms like return rate, Sharpe ratio). For understanding better, I have also checked the difference in meaning between some similar terms MLP and CNN, the former one is fully connected while the latter one is sparsely or partially connected.

More specific to the paper that I have selected, I have learnt that there is existing model which can involve both stock prices and stock news for stock price prediction. I have learnt that the use of dual-channel fusion can avoid interference between two types of features and at the same time gain the benefit from the fusion of two different models.

I have also learnt the text data extraction by using CNN which involves segmentation of text, stopping of words, embedding of words, compression of word vectors into sentence vectors, that of sentence vectors into chapter vectors, and that of chapter vectors into enterprise vector.

Through the stock prediction and the trading strategies parts, I have applied what I have learnt in the lectures about how to calculate various types of technical indicators (e.g. moving averages, RSI, MACD, Bollinger bands, etc.) and how to preliminarily look for buy/sell signals regardless of my limited skill for further checking for uptrend or downtrend from a more macroscopic perspective.

I have also learnt how to carry out mock trading by using Excel.

I have learnt that when copying the data from the table in Powerpoint slides to Excel, I need to need to remove last character in order to have numbers being recognizable by the Excel in order to plot graphs like clustered column bar charts.

I have also learnt that the axes of different charts should be scaled to be the same for fair comparison, and the minimum value of y-axis should be tuned for obvious contrast.

Through the preparation of all reports and materials for this project, I have learnt to cooperate, and share works with my groupmates. Without my groupmates, I think I cannot finish this project by myself. We shared the work well as we did what we are good at. They have trained and tested the models with different features prepared by using the Orange, prepared the Powerpoint slides for presentation and the written report. I have learnt from how they worked on their works as well.

**Individual Report – 21014947G Yu Hin Chung Nikko**

I have done the design of the model. The model is a direct inspiration from the laboratory exercises. The figure is as shown. Some parameters in the CSV file report are inspired by the research paper using ANN and random forest such as moving average. The actual implementation is done by my teammates.
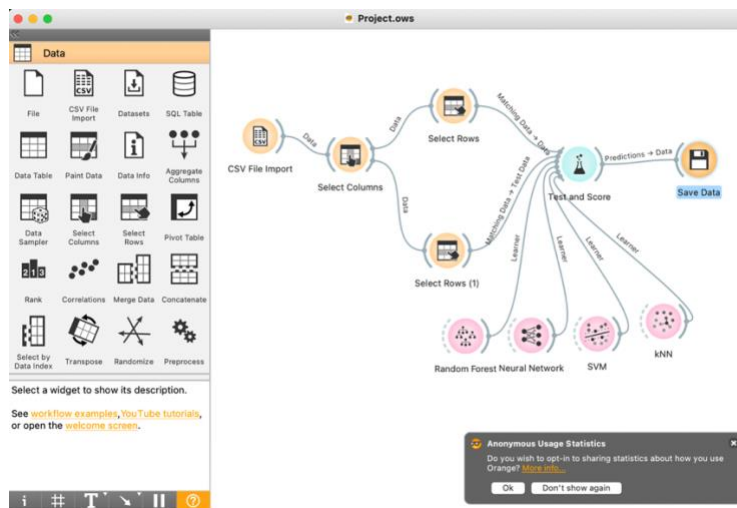


Fig. 31 The model

Meanwhile, I have also assisted the design of the features set, such as suggesting feature set C. Putting the csv files and generating the errors and metrics for feature sets A-D.



Fig. 32 One of my work in generating the result csv and performing MAE metrics for model C.

Apart from my individual work, I have also helped verifying other teammate's model (csvs) such as the curves, the formulas and more. Teammates and I had suggested different strategies during the brainstorming, and we agreed upon making strategies 1-5B.