



Faculty of Computer Science and Information Technology

Predicting US Stock Prices Using Long Short-Term Memory (LSTM)

DAYANG AFIQAH LIYANA BINTI ABANG EHSAN

**Bachelor of Computer Science with Honours
(Multimedia Computing)**

2023

UNIVERSITI MALAYSIA SARAWAK

THESIS STATUS ENDORSEMENT FORM

TITLE Predicting US Stock Prices Using Long Short-Term Memory (LSTM)

ACADEMIC SESSION: 2022/2023

DAYANG AFIQAH LIYANA BINTI ABANG EHSAN

(CAPITAL LETTERS)

hereby agree that this Thesis* shall be kept at the Centre for Academic Information Services, Universiti Malaysia Sarawak, subject to the following terms and conditions:

1. The Thesis is solely owned by Universiti Malaysia Sarawak
2. The Centre for Academic Information Services is given full rights to produce copies for educational purposes only
3. The Centre for Academic Information Services is given full rights to do digitization in order to develop local content database
4. The Centre for Academic Information Services is given full rights to produce copies of this Thesis as part of its exchange item program between Higher Learning Institutions [or for the purpose of interlibrary loan between HLI]
5. ** Please tick (✓)

CONFIDENTIAL (Contains classified information bounded by the OFFICIAL SECRETS ACT 1972)

RESTRICTED (Contains restricted information as dictated by the body or organization where the research was conducted)

UNRESTRICTED



(AUTHOR'S SIGNATURE)

Validated by



DR SHAPIEE ABD RAHMAN
Senior Lecturer (DS2)
Computational Science Programme
Faculty of Computer Science & Information Technology
Universiti Malaysia Sarawak

(SUPERVISOR'S SIGNATURE)

Permanent Address

No.20, Lorong Depo 11, Taman
Sebah Heights, Jalan Depo,
Petra Jaya, 93050, Kuching, Sarawak

Date: 1/7/2023

Date:1/7/2023

Note * Thesis refers to PhD, Master, and Bachelor Degree

** For Confidential or Restricted materials, please attach relevant documents from relevant organizations / authorities

Declaration

I, Dayang Afiqah Liyana Binti Abang Ehsan, hereby declare that the thesis entitled “Predicting US Stock Prices Using Long Short-Term Memory (LSTM)” is entirely based on my own original work, except for appropriate quotations and citations. I also declare that it has not been submitted for any other degree at the Universiti Malaysia Sarawak (UNIMAS) previously or currently.

Signature,



.....

(DAYANG AFIQAH LIYANA BINTI ABANG EHSAN)

Faculty of Computer Science and Information Technology

Date: 25th January 2023

Universiti Malaysia Sarawak

Acknowledgement

First and foremost, all praise to Allah for all His blessings and guidance to me for successfully completing this project after spending about four months of hard works along with numerous challenges and problems that I have encountered.

I would like to express my heartfelt appreciation to my supervisor, Dr. Shapi-ee bin Abd Rahman, and my co-supervisor, Dr. Nurul Syuhada binti Ismail, for granting me the opportunity to conduct this research and for their invaluable guidance and support throughout the process. I am deeply grateful to both Dr. Shapi-ee bin Abd Rahman and Dr. Nurul Syuhada binti Ismail for being a constant source of inspiration and motivation, and for providing valuable references that have contributed to the completion of this research.

My deepest gratitude for my parents and my family, Izzati Nadhirah Binti Mohammad Saberi, Anis Hazirah Binti Mohammad Saberi, Yusha 'Athirah Binti Rozman and Husna Amani Binti Rozman. Special thanks to my friends, Nur'azra Alia Nisa Binti Zulpakar, Syahazwani Nurain Binti Shariman Faizul, Nurfazlina Binti Yusuf, Nurin Alya Binti Haris, Siti Rubiah Binti Muslim, Dayang Nurazzyati Binti Awang Suffian and Elia Sari Binti Odita who also have helped me a lot in giving guidance and information regarding my research.

Lastly thank you to those who are not mentioned above and whoever was involved directly and indirectly in completing this final year project.

Thank You,

DAYANG AFIQAH LIYANA BINTI ABANG EHSAN

Abstract

A stock, commonly referred to as equity, represents an investment that signifies partial ownership in a company. Investors are concerned with two crucial aspects: the current price of their existing or potential investment and its projected selling price in the future. Predicting stock prices has always been of great interest to investors; however, it has proven to be a challenging task for researchers and analysts. The stock market is highly unpredictable, with numerous complex financial indicators. Consequently, financial analysts, researchers, and data scientists are continuously exploring analytical tools to uncover stock market patterns. In this study, historical stock price data is leveraged to predict the stock prices of selected US companies using a machine learning approach known as the Long Short-Term Memory (LSTM) Model, which is a specialized form of Recurrent Neural Network (RNN). The dataset comprises five years of AAPL and MSFT data obtained from Yahoo Finance, with consideration given to six relevant attributes. The LSTM model is employed to generate accurate and reliable predictions. The LSTM model holds several advantages in the realm of stock price prediction as it utilises historical stock price data to discern patterns and trends, enabling the forecasting of future price movements. This study focuses on employing the LSTM model to shed light on the potential for achieving precise stock price forecasts using machine learning techniques. Additionally, the Root Mean Square Error (RMSE) is employed as a supplementary performance measure alongside the LSTM model for stock price prediction.

Abstrak

Saham, biasanya dirujuk sebagai ekuiti, mewakili pelaburan yang menandakan pemilikan separa dalam syarikat. Pelabur bimbang dengan dua aspek penting: harga semasa pelaburan sedia ada atau berpotensi mereka dan harga jualan yang diunjurkan pada masa akan datang. Meramalkan harga saham sentiasa menarik minat pelabur; Walau bagaimanapun, ia telah terbukti menjadi tugas yang mencabar bagi penyelidik dan penganalisis. Pasaran saham sangat tidak dapat diramalkan, dengan banyak petunjuk kewangan yang kompleks. Akibatnya, penganalisis kewangan, penyelidik, dan saintis data terus meneroka alat analisis untuk mendedahkan corak pasaran saham. Dalam kajian ini, data harga saham sejarah dimanfaatkan untuk meramalkan harga saham syarikat AS terpilih menggunakan pendekatan pembelajaran mesin yang dikenali sebagai Model Memori Jangka Pendek Panjang (LSTM), yang merupakan bentuk khusus Rangkaian Neural Berulang (RNN). Set data terdiri daripada lima tahun data AAPL dan MSFT yang diperoleh daripada Yahoo Finance, dengan pertimbangan diberikan kepada enam atribut yang berkaitan. Model LSTM digunakan untuk menghasilkan ramalan yang tepat dan boleh dipercayai. Model LSTM memegang beberapa kelebihan dalam bidang ramalan harga saham kerana ia menggunakan data harga saham sejarah untuk membezakan corak dan trend, membolehkan ramalan pergerakan harga masa depan. Kajian ini memberi tumpuan kepada penggunaan model LSTM untuk memberi penerangan tentang potensi untuk mencapai ramalan harga saham yang tepat menggunakan teknik pembelajaran mesin. Selain itu, Ralat Root Mean Square (RMSE) digunakan sebagai langkah prestasi tambahan bersama model LSTM untuk ramalan harga saham.

Table of Contents

Table of Contents

CHAPTER 1: INTRODUCTION	1
1.1 INTRODUCTION	1
1.2 PROBLEM STATEMENT	2
1.3 SCOPE	5
1.4 OBJECTIVES	5
1.5 BRIEF METHODOLOGY	6
1.6 SIGNIFICANCE OF PROJECT	8
1.7 PROJECT SCHEDULE	10
1.8 EXPECTED OUTCOME	12
1.9 PROJECT SUMMARY	12
CHAPTER 2: LITERATURE REVIEW	13
2.1 INTRODUCTION	13
2.2 STOCK AND STOCK PRICES	13
2.3 MACHINE LEARNING	20
<i>2.3.1 Types of Machine Learning</i>	<i>22</i>
<i>2.3.2 Types of Algorithms that used Supervised and Unsupervised Machine Learning</i>	<i>24</i>
<i>2.3.3 Using Machine Learning to Predict Stock Prices</i>	<i>25</i>
<i>2.3.4 Recurrent Neural Network (RNN)</i>	<i>26</i>
<i>2.3.5 Long Short-Term Memory (LSTM) Model</i>	<i>29</i>
2.4 REVIEW OF SIMILAR RESEARCH	31
<i>2.4.1 Predicting Stock Prices Using Artificial Neural Network (ANN)</i>	<i>32</i>
<i>2.4.2 Predicting Stock Prices using Random Forest</i>	<i>39</i>
<i>2.4.3 Predicting Stock Prices using Long Short-Term Memory (LSTM) Model</i>	<i>46</i>
2.5 COMPARISON OF RELATED RESEARCH WORK	51
2.6 CONCLUSION	53
CHAPTER 3: METHODOLOGY	54
3.1 INTRODUCTION	54
3.2 SOURCES AND SOFTWARE	54
3.3 LONG SHORT-TERM MEMORY (LSTM) ARCHITECTURE	55
3.4 FRAMEWORK OF STUDY	57
3.5 PARAMETERS USED	58
3.6 DATA COLLECTION	58
3.7 DATA PRE-PROCESSING	60
3.8 MODELLING LONG SHORT-TERM MEMORY (LSTM) MODEL	62
3.9 TRAINING THE LONG SHORT-TERM MEMORY (LSTM) MODEL	64
3.10 MODEL EVALUATION AND VISUALISATION	65
3.11 COMPARISON BETWEEN MODELS	67
3.12 CONCLUSION	67
CHAPTER 4: IMPLEMENTATION	68

4.1 INTRODUCTION	68
4.2 DATA COLLECTION	68
4.3 DATA PRE-PROCESSING	70
4.3.1 Data Cleaning	70
4.3.2 Importing cleaned dataset into RapidMiner	75
4.3.3 Set Roles of Attributes	78
4.3.4 Splitting the dataset in training and testing datasets	80
4.3.5 Featuring selection	81
4.4 CROSS VALIDATION PROCESS	82
4.4.1 Modelling the LSTM model	85
4.4.2 Modelling the LSTM layer	92
4.4.3 Predictive LSTM Model	100
4.4.4 Apply Model	103
4.4.5 Performance Measure	105
4.5 CONCLUSION	106
CHAPTER 5: RESULT AND DISCUSSION	107
5.1 INTRODUCTION	107
5.2 RESULT AND DISCUSSION	108
5.3 CONCLUSION	166
CHAPTER 6: CONCLUSION	168
6.1 INTRODUCTION	168
6.2 CONCLUSION	169
6.3 LIMITATIONS	170
6.4 FUTURE WORKS	173
REFERENCES	176

List of Tables

Table 2.1: Types of Algorithms in Supervised and Unsupervised Learning.....	24
Table 2.2: Researchers' LSTM Model Summary (Moghar & Hamiche, 2020).....	47
Table 2.3: The value of loss for GOOGL and NKE for different number of epochs (Moghar & Hamiche, 2020).....	50
Table 2.4: Comparison of related research work.....	51
Table 3.1: List of parameters used.....	58
Table 3.2: Data pre-processing process.....	60
Table 5.1: Predicted close price and the actual close price of Apple Inc for all the dates (including training data and testing data, excluding weekends and public holidays).....	113
Table 5.2: Predicted close price and the actual close price of Apple Inc for all the dates (including training data and testing data, excluding weekends and public holidays) when RMSE value of \$1.113 is obtained.....	123
Table 5.3: Predicted close price and the actual close price of Apple Inc for all the dates (including training data and testing data, excluding weekends and public holidays) when RMSE value of \$0.986 is obtained.....	128
Table 5.4: Predicted close price and the actual close price of Apple Inc for all the dates (including training data and testing data, excluding weekends and public holidays) when RMSE value of \$1.008 is obtained.....	133
Table 5.5: Predicted close price and the actual close price of Microsoft Corporation for all the dates (including training and testing data, excluding weekends and public holidays).....	141
Table 5.6: Predicted close price and the actual close price of Microsoft Corporation for all the dates (including training and testing data, excluding weekends and public holidays) when RMSE value of \$1.963 is obtained.....	150
Table 5.7: Predicted close price and the actual close price of Microsoft Corporation for all the dates (including training and testing data, excluding weekends and public holidays) when RMSE value of \$2.034 is obtained.....	155
Table 5.8: Predicted close price and the actual close price of Microsoft Corporation for all the dates (including training and testing data, excluding weekends and public holidays) when RMSE value of \$1.884 is obtained.....	160

List of Figures

Figure 1.1: Framework of the study, starts from data collection to data pre-processing to splitting data and then modelling and training the LSTM model and lastly to model visualisation and performance evaluation.....	6
Figure 1.2: Gantt Chart for FYP 1	10
Figure 1.3: Gantt Chart for FYP2	11
Figure 2.1: E*TRADE signup page (E*TRADE, 2023).....	16
Figure 2.2: Firstrade homepage (Firstrade, 2023)	16
Figure 2.3: Charles Schwab homepage (Charles Schwab, 2023)	17
Figure 2.4: Robinhood homepage (Robinhood, 2023)	17
Figure 2.5: RNN Forward Propagation Structure (Du et al., 2019).....	27
Figure 2.6: RNN Forward Propagation Timing Diagram (Du et al., 2019).....	28
Figure 2.7: Internal Structure of LSTM (Du et al., 2019).....	29
Figure 2.8: Training phase of the researchers' model (Haider Khan et al., 2011).....	34
Figure 2.9: Single neuron used in the researchers' model (Haider Khan et al., 2011).....	35
Figure 2.10: Graphical representation of Predicting and Actual price of ACI pharmaceutical using 2 input data sets to predict the stock values for first 8 days of November 2010 where the input past historical data is from 31-08-2010 to 30-09-2010 (Haider Khan et al., 2011)	36
Figure 2.11: Graphical representation of Predicting and Actual price of ACI pharmaceutical using 5 input data sets to predict the stock values for first 8 days of November 2010 where the input past historical data is from 1-9-2010 to 31-10-2010 (Haider Khan et al., 2011)	37
Figure 2.12: Random Forest Classifier Algorithm (Khaidem et al., 2016)	43
Figure 2.13: OOB error rate vs Number of estimators (Khaidem et al., 2016)	43
Figure 2.14: Results from Dai and Zhang (2013) (Khaidem et al., 2016).....	44
Figure 2.15: Results for 3M stock obtained with the researchers' model (Khaidem et al., 2016)	45
Figure 2.16: Researchers' LSTM Model Structure (Moghar & Hamiche, 2020).....	48
Figure 2.17: Result of training for the NKE stocks with different dataset time (Moghar & Hamiche, 2020).....	48
Figure 2.18: Result of training for NKE and GOOGL stocks with different number of epochs (Moghar & Hamiche, 2020).....	49
Figure 3.1: Repeating Module in RNN reliance (Istiake Sunny et al., 2020).....	55
Figure 3.2: Repeating Module in LSTM (Istiake Sunny et al., 2020)	56
Figure 3.3: Framework of the study, starts from data collection to data pre-processing to splitting data and then modelling and training the LSTM model and lastly to model visualisation and performance evaluation.....	57
Figure 3.4: Apple Inc. (AAPL) stock price in Yahoo Finance (Yahoo Finance, 2023)	59

Figure 3.5: Microsoft Corporation (MSFT) stock price in Yahoo Finance (Yahoo Finance, 2023)	59
Figure 4.1: Historical stock price for AAPL.....	69
Figure 4.2: Historical stock price for MSFT.....	69
Figure 4.3: Missing values for date column for Apple Inc. and Microsoft Corporation dataset	71
Figure 4.4: Dates column for Apple Inc. and Microsoft Corporation dataset after replacing the missing value	71
Figure 4.5: Apple Inc. dataset before removing unnecessary columns	72
Figure 4.6: Microsoft Corporation dataset before removing unnecessary columns	72
Figure 4.7: Apple Inc. dataset after removing unnecessary columns	73
Figure 4.8: Microsoft Corporation . dataset after removing unnecessary columns.....	73
Figure 4.9: Previous close price column for Apple Inc. dataset	74
Figure 4.10: Previous close price column for Microsoft Corporation dataset.....	75
Figure 4.11: Read Excel operator	76
Figure 4.12: Importing the Apple.Inc dataset	76
Figure 4.13: Format column step for Apple Inc. dataset	77
Figure 4.14: Connection between the ‘Read Excel’ operator and ‘Set Role’ operator’	78
Figure 4.15: Assigning roles to the attributes	79
Figure 4.16: Connection of the ‘Set Role’ operator and ‘Split Data’ operator.....	80
Figure 4.17: Splitting the data into 70% for training and 30% for testing.....	81
Figure 4.18: Connection of ‘Select Attributes’ operator and ‘Split Data’ operator.....	82
Figure 4.19: Parameters configuration for feature selection process.....	82
Figure 4.20: Connection between the ‘Select Attributes’ and ‘Cross Validation’ operator	84
Figure 4.21: Parameters for the cross validation operator	84
Figure 4.22: Connection of the ‘Deep Learning’ operator inside the ‘Cross Validation’ operator	91
Figure 4.23: Parameters that has been configured for the ‘Deep Learning’ operator.....	92
Figure 4.24: Connection of the ‘Add LSTM Layer’ operator and the ‘Add Output Layer’ operator inside the ‘Deep Learning’ operator	98
Figure 4.25: Parameters that has been tuned for the ‘Add LSTM Layer’ operator	99
Figure 4.26: Parameters that has been tuned for the ‘Add Output Layer’ operator.....	99
Figure 4.27: Connection of the ‘Deep Learning’ operator to the ‘Predictive Deep Learning’ operator inside the ‘Cross Validation’ operator	102
Figure 4.28: Parameters that has been tuned for the ‘Predictive Deep Learning’ operator...	103
Figure 4.29: Connection of the ‘Apply Model’ operator inside the ‘Cross Validation’ operator	104

Figure 4.30: Connection of the ‘Performance’ operator inside the ‘Cross Validation’ operator	106
Figure 5.1: RMSE value for Apple.Inc dataset.....	108
Figure 5.2: Predicted close price and the actual close price of Apple Inc.....	112
Figure 5.3: Graph of predicted close price and actual price of Apple Inc.	118
Figure 5.4: RMSE value of \$1.113 obtained on different execution of the project.....	121
Figure 5.5: RMSE value of \$0.986 obtained on different execution of the project.....	121
Figure 5.6:RMSE value of \$1.008 obtained on different execution of the project.....	121
Figure 5.7: Predicted close price and the actual close price of Apple Inc when RMSE value of \$1.113 is obtained	122
Figure 5.8: Predicted close price and the actual close price of Apple Inc when RMSE value of \$0.986 is obtained	127
Figure 5.9: Predicted close price and the actual close price of Apple Inc when RMSE value of \$1.008 is obtained	132
Figure 5.10: RMSE value obtains for Microsoft Corporation dataset.....	138
Figure 5.11: Predicted close price and the actual close price of Microsoft Corporation when a RMSE value of \$1.865 is obtained	140
Figure 5.12: Graph of predicted close price and actual price of Microsoft Corporation when a RMSE value of \$1.865 is obtained	145
Figure 5.13: RMSE value of \$1.963 obtained on different execution of the project.....	148
Figure 5.14: : RMSE value of \$2.034 obtained on different execution of the project.....	148
Figure 5.15: : RMSE value of \$1.884 obtained on different execution of the project.....	148
Figure 5.16: Predicted close price and the actual close price of Microsoft Corporation when RMSE value of \$1.963 is obtained	149
Figure 5.17: Predicted close price and the actual close price of Microsoft Corporation when RMSE value of \$2.034 is obtained	154
Figure 5.18: Predicted close price and the actual close price of Microsoft Corporation when RMSE value of \$1.884 is obtained	159

List of Equations

Equation 2.1	30
Equation 2.2	30
Equation 2.3	30
Equation 2.4	30
Equation 2.5	31
Equation 2.6	31
Equation 2.7	33
Equation 2.8	35
Equation 2.9	35
Equation 2.10	40
Equation 2.11	40
Equation 2.12	41
Equation 2.13	41
Equation 2.14	42
Equation 2.15	42
Equation 3.1	56
Equation 3.2	57
Equation 3.3	57
Equation 3.4	57
Equation 3.5	66
Equation 4.1	97
Equation 5.1	110

Chapter 1: Introduction

1.1 Introduction

Every investor should consider two crucial aspects: the current price of their investment or the one they plan to buy, and its future selling price. However, investors consistently examine past pricing history and use it to influence their future investment decisions. Some investors avoid purchasing stocks or indices that have experienced rapid price increases because they expect a correction. Conversely, others avoid buying stocks that are in decline, assuming the downward trend will persist. The primary goal of stock price prediction is to generate substantial profits (Hagenau et al., 2013). Predicting the performance of the stock market poses a challenging task. Various factors, including physical and psychological influences, rational and irrational behaviour, and others, also affect these predictions. All these factors contribute to the dynamic and volatile nature of stock prices. Consequently, accurately predicting stock prices becomes extremely difficult (Yu & Yan, 2019). Due to the abundance of data available, there is ongoing exploration by financial analysts, researchers, and data scientists to uncover patterns in the stock market. The potential of data analysis to be a game changer in this field has been investigated. Numerous studies have focused on utilising machine learning in quantitative finance, which allows for predicting the prices of various assets, optimizing investment strategies, and other related processes (Moghar & Hamiche, 2020). Generally, machine learning refers to algorithmic techniques that employ computers to discover patterns solely based on data, without requiring explicit programming instructions. These techniques have the capability to unveil previously unnoticed patterns and insights, thereby enabling the generation of highly accurate forecasts (Ghosh et al., 2019). In the field of quantitative finance, there are numerous models available that can be combined with machine learning to predict the future value of assets. These models offer a way to integrate diverse sources of information, creating a unique and effective tool. Recent advancements in

machine learning techniques, such as neural networks, gradient boosted regression trees, support vector machines, and random forest, have emerged from the fusion of statistics and learning models. These algorithms are capable of capturing complex patterns characterized by non-linearity and uncovering relationships that are challenging to identify using linear algorithms. Moreover, these algorithms demonstrate superior effectiveness and handle multicollinearity better than linear regression techniques. The research paper "Trends and Applications of Machine Learning in Quantitative Finance" (Emerson et al., 2019) focuses primarily on the investment process, which is a fundamental component of finance. This encompasses predicting investment returns, modelling risks, and constructing portfolios. In this regard, a group of machine learning algorithms based on Recurrent Neural Network (RNN) proves to be valuable for forecasting and predicting financial market prices. Another study examines the accuracy of autoregressive integrated moving average (ARIMA) and long short-term memory (LSTM) as prediction algorithms for time series data. When applied to financial data, it was found that LSTM outperformed ARIMA by a significant margin (Siami Namin & Siami Namin, 2018). The objective of this project is to utilise a machine learning algorithm based on LSTM RNN to forecast stock prices in the US stock market and assess the performance of the Long Short-Term Memory (LSTM) Model.

1.2 Problem Statement

Investors have a keen interest in predicting stock prices as they seek profitable investment opportunities. However, many investors lack a reliable strategy to make informed decisions that yield significant profits. Consequently, there is a growing interest among investors to gain insights into the future state of the stock market, enabling them to make intelligent and lucrative investments. Nevertheless, predicting stock prices has proven to be a challenging task for researchers and analysts. The stock market is influenced by a multitude of

factors, including both psychological and physical elements, rational and irrational behaviour, investor sentiment, and more. The unpredictable nature of stock market fluctuations, coupled with the complexity of numerous financial indicators, further compounds the difficulty of accurate prediction. Consequently, there is a need to analyse vast amounts of data to identify patterns in the stock market. Financial analysts, researchers, and data scientists continue to explore analytics tools in their quest to uncover stock market patterns. The application of data analysis in this domain has the potential to bring about significant advancements. To predict future trends, researchers employ Machine Learning techniques using historical stock price data. Leveraging these Machine Learning approaches allows for the discovery of previously unnoticed patterns and insights, enabling the generation of highly accurate forecasts (Ghosh et al., 2019). For this project, the Long Short-Term Memory (LSTM) Model will be employed to predict stock prices. LSTM is a specialized variant of recurrent neural network (RNN) architecture designed to better capture temporal sequences and their long-term relationships compared to traditional RNNs (Van Houdt et al., 2020; Sak et al., 2014). RNNs, in general, are neural networks devised to overcome the limitations of conventional neural networks. They possess a "memory" that stores information about previous calculations. In essence, RNNs are a type of artificial neural network in which the connections between units form a directed graph along a sequence (Patel et al., 2018). These networks incorporate internal feedback loops that facilitate information retention. Theoretically, RNNs are capable of handling "long-term dependencies" effectively. Unfortunately, in practice, RNNs seem unable to effectively learn long-term dependencies, giving rise to what is known as the "Vanishing gradient problem." The neural network employs the gradient descent algorithm to update its weights. However, as the network progresses through lower layers, the gradients diminish significantly. Consequently, the network reaches a point where there is little room for improvement as the gradients remain almost constant. This adjustment adversely affects the network's output.

When the gradient difference becomes exceedingly small, the network fails to learn anything new, and the output remains unchanged. On the other hand, exploding gradients lead to unstable training due to fluctuating weights, while vanishing gradients result in excessively long training times and, in some cases, render training impossible (Tsun et al., 2019). Consequently, a network grappling with vanishing gradient problems fails to converge on an effective solution, rendering both scenarios undesirable in neural network training. Hence, it necessitates the development of novel training methods and architectures to tackle these challenges. In response, in 1997, Hochreiter and Schmidhuber introduced the Long Short-Term Memory (LSTM) Model, a specialized type of RNN that overcomes the limitations of traditional RNNs by effectively learning long-term dependencies. Over time, LSTMs have undergone refinements to address the challenge of handling long-term dependencies. The LSTM architecture tackles this issue by incorporating memory cells and gate units into neural networks (Tsun et al., 2019). LSTMs were specifically designed with a unique "memory cell" structure that enables the retention or discarding of information as it traverses through the network. Furthermore, LSTMs employ three distinct gates (input, forget, and output) to regulate the flow of data and determine what information to remember and what to forget. In a typical formulation, memory cells store information encountered in their cell state. When an input is received, the output is determined based on a combination of the cell state (representing previous information) and the updated cell state. As subsequent inputs are fed into the memory cell, the updated state and new input can be utilised to compute the updated output. With the ability to retain long-term dependencies, LSTMs are well-suited for tasks like time series prediction, such as forecasting a company's share price. Unlike traditional RNNs, LSTMs have gate mechanisms that govern the transfer of data, providing finer control over what information is stored and what is disregarded. This enhanced control allows LSTMs to effectively handle non-linear relationships present in financial data, thereby enabling more accurate predictions.

Moreover, LSTM models demonstrate proficiency in handling extensive datasets with numerous time-steps, leveraging historical data to make more precise forecasts regarding future trends. Consequently, LSTM emerges as a favourable choice for predicting the future performance and growth of a company's share price.

1.3 Scope

This project's scope is to model, tune and evaluate the performance of Long Short-Term Memory (LSTM) Model, which is a special type of Recurrent Neural Network, in predicting the stock prices of companies listed on the US stock markets based on data on the date, opening stock price, closing stock price, previous close price, lowest price and highest price. This project will use five years of data in the prediction phase. Some of the software and sources that will be used to conduct this project includes Yahoo Finance and RapidMiner.

1.4 Objectives

The objectives of this project are:

- To collect data of several public companies listed on the US stock exchange.
- To predict the stock prices of the selected companies using Long Short-Term Memory (LSTM).
- To evaluate the performance of the Long Short-Term Memory (LSTM) Model.

1.5 Brief Methodology

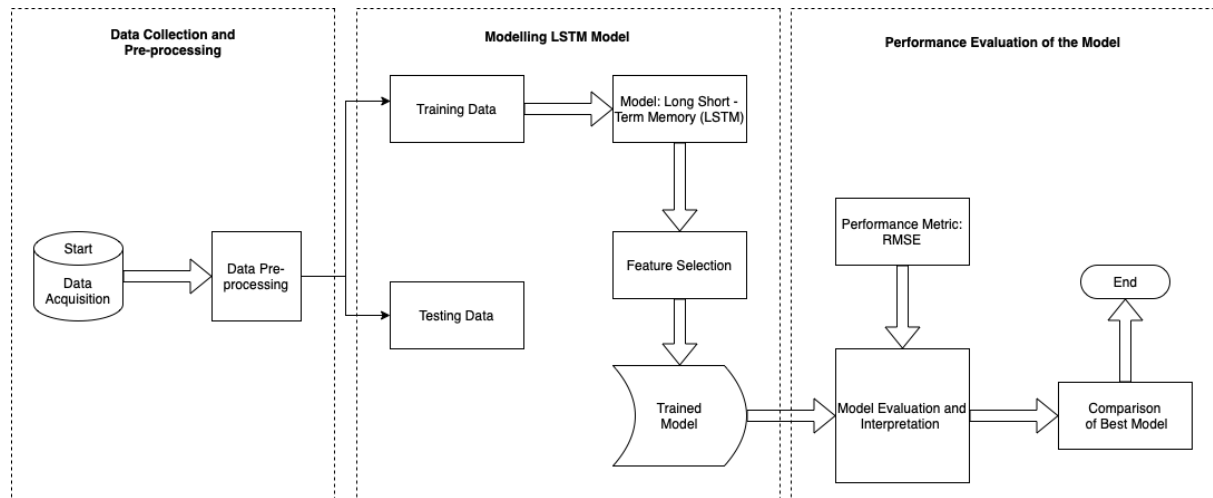


Figure 1.1: Framework of the study, starts from data collection to data pre-processing to splitting data and then modelling and training the LSTM model and lastly to model visualisation and performance evaluation

Phase 1: Data Collection and Pre-processing

For this project, the first phase is to collect 5 years' worth of data of past stock prices from selected companies using Yahoo Finance. Date, open share price, close share price, high share price, low share price and previous close prices are used.

The next part is data pre-processing. Before modelling the LSTM model, it is important to conduct data pre-processing. It is a crucial stage in machine learning research. Frequently, data collection procedures are poorly governed, resulting in out-of-range numbers, missing values, etc. Analysing data that has not been thoroughly examined for such issues can lead to incorrect conclusions. Before conducting an analysis, the representation and quality of the data are therefore of the utmost importance. Data pre-processing is the most crucial element of a machine learning project, particularly when dealing with computational data. Data pre-processing involves data cleaning, splitting data into training and testing set and featuring selection.

Phase 2: Modelling and Training

Moving to the subsequent phase involves the modelling and training of the LSTM model. The LSTM units are employed as foundational components for the layers of a recurrent neural network (RNN), commonly known as an LSTM network. These LSTM units enable RNNs to retain inputs over an extended period by leveraging a computer-like memory structure. Within an LSTM unit, information can be written, deleted, and read from this memory. Moreover, the LSTM memory cell encompasses three gates: the forget gate, input gate, and output gate. The forget gate determines when specific elements of the cell state should be replaced with more recent data. The input gate governs the conditions for storing or updating data in the cell state. Finally, the output gate determines the information transmitted to the subsequent network node (Nivethitha et al., 2019). In featuring selection process earlier, data attributes are chosen and are going to be fed to the LSTM model. Since this project utilise RapidMiner, there is no need to create the LSTM model from scratch since the platform offers pre-built LSTM model components that are specifically designed for machine learning applications like stock price prediction. Therefore, in this phase it would be more to configuring, optimizing, and training the existing LSTM model in RapidMiner. Modelling involves customizing the architecture and parameters of the LSTM model to suit the specific requirements of the predicting task and data. This includes making choices regarding the number of layers, neurons, activation functions, and other architectural elements. After configuring the model, the next step is tuning, which aims to optimize the model's performance. Tuning entails experimenting with various hyperparameter settings, such as learning rate, batch size, regularization techniques, and optimization algorithms. The goal is to identify the optimal combination of hyperparameters that generates the most accurate predictions. Once the modelling and tuning stages are complete, the LSTM model is ready for training. During training, the model is fed with the training data, enabling it to learn and adjust its internal

parameters to capture the underlying patterns and trends in the historical stock price data. This training process equips the model with the ability to make predictions based on the acquired patterns.

Phase 3: Performance Evaluation and Model Visualisation

In this phase, the trained LSTM model will be evaluated with the test set and the Root Mean Square Error (RMSE) metric is used to assess the model's performance. RMSE is a commonly used metric for evaluating the accuracy of a model's predictions. It measures the differences between actual and predicted values, known as residuals. Usually, it is beneficial to graphically depict the predicted prices. Stock price data will be represented as a 2D graph in RapidMiner where the current and predicted close price data will be compared. Predictions and existing data will be utilised to determine precision.

1.6 Significance of Project

This research on stock market prediction offers the potential advantage of significantly reducing losses. Many investors make the mistake of insufficient research before attempting to predict the market, leading them to use ineffective forecasting methods. They often rely on gut feelings or arbitrary estimates when investing in stocks, hoping for price increases and profits. However, by understanding and utilising appropriate prediction strategies, investors can minimize their losses. Making well-informed decisions using reliable stock market prediction tools has the potential to substantially increase profits.

Furthermore, stock market forecasting provides the benefit of consistent results. Given the high volatility of the stock market, there is no guarantee that investors will always be profitable, even if they employ various tactics and formulas to forecast the future. They may experience gains on some days and losses of equal magnitude on others. In such situations, consistency becomes crucial, regardless of the extent of gains and losses, in order to

consistently provide investors with favourable returns. Accurate predictions of the stock market help investors achieve the consistency necessary to generate higher returns in this volatile market.

1.7 Project Schedule

Figure 1.2 and 1.3 depicts the Gantt Chart for this project:

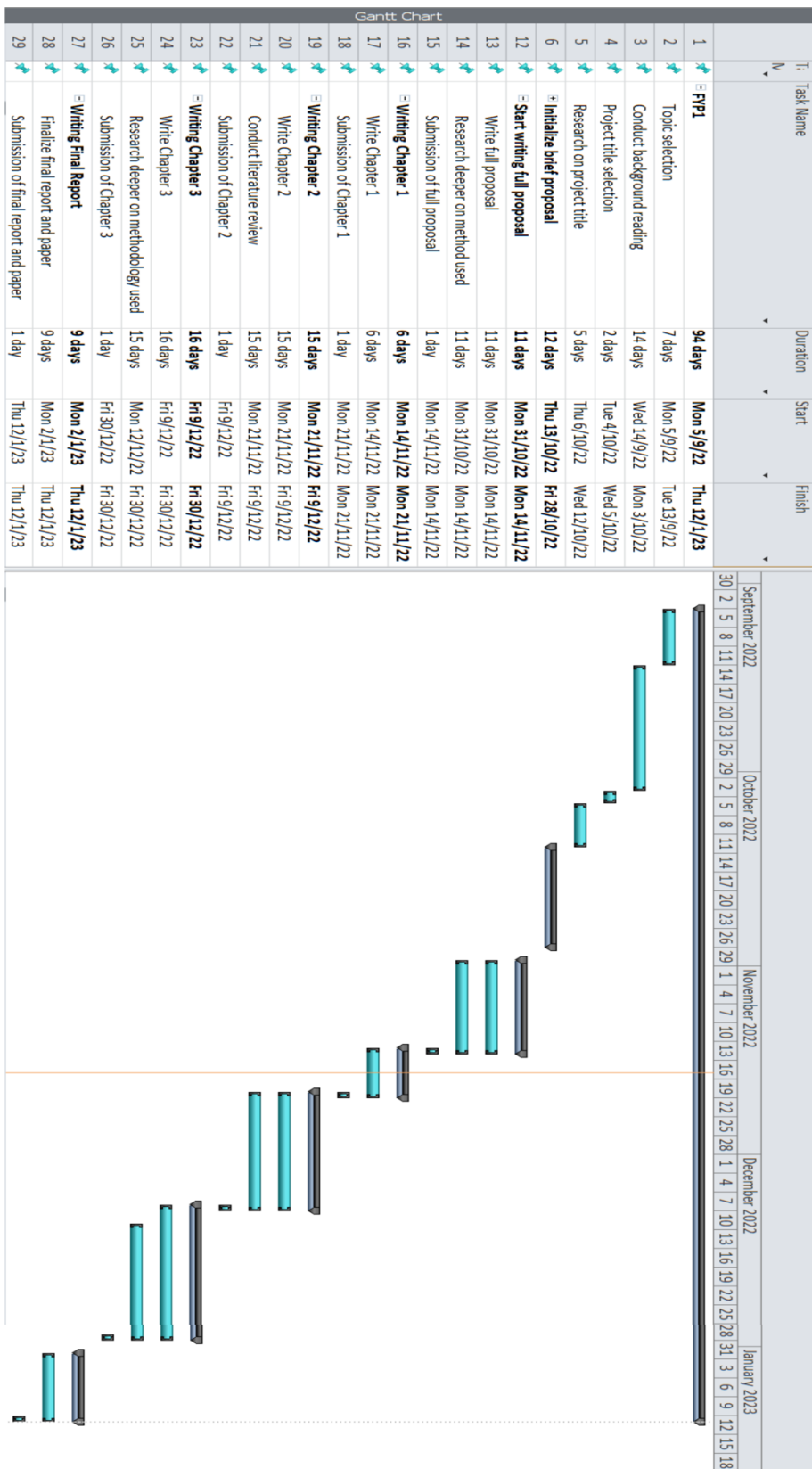


Figure 1.2: Gantt Chart for FYP 1

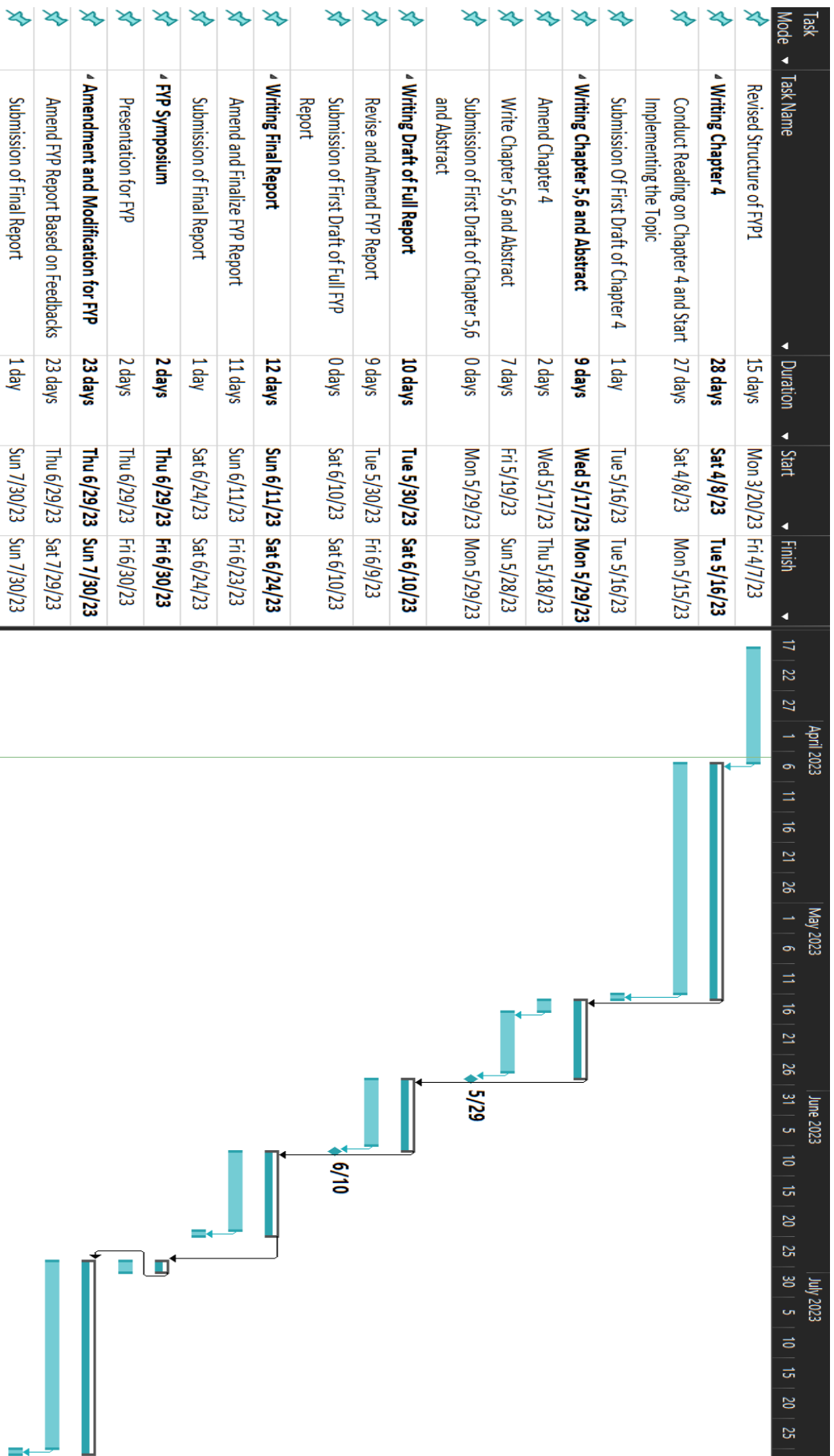


Figure 1.3: Gantt Chart for FYP2