

**FINANCIAL NUMERIC AND TEXTUAL DATA BASED  
STOCK PREDICTION USING MACHINE LEARNING  
TECHNIQUES**

**BY**

**MOHAMMAD RABIUL ISLAM**

A thesis submitted in fulfilment of the requirement for the  
degree of Doctor of Philosophy in Computer Science

Kulliyyah of Information and Communication Technology  
International Islamic University Malaysia

December 2020

## ABSTRACT

Source of web-text and numerical data analysis for stock prediction is the challenging tasks in today's stock market engineering. Day traders in the stock market face the common issues of decision making, which is mostly dependent on daily or weekly bases data analysis. To overcome this problem, web text mining and data mining analysis techniques applied on stock market closing values, which brings the most technical approach. In terms of stock-market textual data classification, the applicable soft-computing technique holds two classifications as binary and multinomial with clustering algorithms used to apply for analysis. New prediction models overcome the drawback of previous research and indicate the necessity of classification by creating prediction algorithm with a token or a polar based financial text weighting scheme of intensive scale value (ISV) system. Binary classification helps to improve the sense of positivity and negativity with intensive value to evaluate the vast amount of financial textual data for trading decision. This research improves with the technical correlation that addressed the problem of categorical financial textual and numerical data throughout various soft-computing techniques. Targeted numerical and textual data rely on subsequently neural network, binary and multinomial classification to improve the prediction techniques by feature engineering. In terms of textual data, the novel financial data dictionary is prepared based on Harvard reference weighting scheme valued that defined as a likely result in new Stock Prediction Model. Financial news-based text analysis techniques improve the classification scenario with Naïve Bayes binary classification through financial data dictionary. Beside the text analysis, the feed-forward neural network architectural model also improved based on backpropagation neural network structural that approached by defining the correlation between the actual and prediction trend of a daily basis. Daily stock price prediction is the main objective of this research and very essential to generate accurate prediction through online daily basis financial news data. The new architectural neural network model performs with sequential data as hidden with the dataset which applied by the multi-objective optimization algorithm. Throughout feature engineering, setting by scaling value determines the weight factors of developing a neural network that used to define more precious trend within this model. This model enabled to calculate the highest frequent value that occurred on a large dataset, and clustering indicate the stock trend of the prediction. Based on the numerical financial data, new Stock Prediction Model (SPM) have developed for analyzing market movement from two benchmarks numerical stock market dataset those are S&P 500index and OHLCV dataset. Developing integrated classification techniques conducting with prediction analysis based on its classification accuracy as defined in this research 82% which is obvious and better than previous research. The performance with feature engineering in text classification also gain 93%, whereas multilevel and binary classification have performed as combined to gain the best accuracy level. Performance of the proposed approach was estimated by evaluating various parameter as part of the information retrieval field, as seen in this experimental result. However, developing model impacts on academical research philosophy in terms of financial data classification but not highly recommend using in real trading analysis.

## خلاصة البحث

تعد المصادر الإلكترونية للنصوص والتحليل الرقمية للبيانات للتنبؤ بالأسهم مهمة صعبة في هندسة سوق الأوراق المالية في وقتنا هذا. حيث يواجه المتداولون اليوميون في سوق الأسهم قضايا مشتركة لاتخاذ القرارات والتي تعتمد في الغالب على تحليل البيانات اليومية أو الأسبوعية لسوق الأوراق المالية. وللتغلب على هذه المشكلة، يتم تطبيق تقنيات التنقيب عن البيانات وتحليلها على قيم إغلاق هذه الاسواق، وهو ما يحقق أفضل نصح تقني. وفيما يتعلق بتصنيف البيانات النصية لسوق الأسهم، فإن الأساليب الحاسوبية المعمول بها تحمل صنفين أحدهما ثنائي والآخر متعدد الحدود مع خوارزميات التجميع المستخدمة للتحليل. وتتغلب نماذج التنبؤ الجديدة على عيوب البحوث السابقة كما تشير إلى ضرورة التصنيف عن طريق إنشاء خوارزمية التنبؤ مع رمز أو نظام ترجيح للنص المالي القائم على نظام قيمة المقياس المكثف (ISV). اما التصنيف الثنائي يساعد على تحسين المعنى الإيجابي والسليبي مع أهمية مكثفة لتقييم الكم الهائل من البيانات النصية المالية لاتخاذ قرار التداول. يتم تحسين هذه الدراسة من خلال الارتباط الفني الذي تناول مشكلة البيانات المالية الرقمية والنصية عبر مختلف التقنيات الحاسوبية. حيث تعتمد البيانات الرقمية والنصية المستهدفة على الشبكات العصبية والتصنيف الثنائي والمتعدد الحدود لتحسين تقنيات التنبؤ من خلال هندسة الميزات. ومن حيث البيانات النصية، تم إعداد قاموس البيانات المالية الجديد استناداً إلى مخطط ترجيح المرجح في جامعة هارفارد، والذي تم تحديده باعتباره نتيجة محتملة لنموذج التنبؤ بالأسهم الجديد. تعمل تحليل بيانات النصوص المالية المستندة على الأخبار المالية على تحسين سيناريو التصنيف باستخدام تصنيف Naïve Bayes الثنائي من خلال قاموس البيانات المالية. وإلى جانب تحليل النص، تم أيضاً تحسين بنية الشبكة العصبية للتغذية المسبقة استناداً إلى بنية الشبكة العصبية العكسية والتي تمت من خلال تحديد العلاقة المتبادلة بين الاتجاه الفعلي والتوقع للأساس اليومي. كما يعد التنبؤ بأسعار الأسهم اليومية الهدف الرئيسي لهذا البحث وهو ضروري للغاية لتقديم تنبؤ دقيق من خلال بيانات الأخبار المالية اليومية على الإنترنت. حيث يعمل نموذج الشبكة العصبية الجديد ببيانات متسلسلة مخفية مع مجموعة البيانات التي يتم تطبيقها بواسطة خوارزمية التحسين متعددة الأهداف. وخلال هندسة الميزات، يتم تحديد الإعدادات حسب قياس قيمة عوامل الوزن لتطوير شبكة عصبية تستخدم لتحديد اتجاه أكثر قيمة ضمن هذا النموذج. تم تمكين هذا النموذج لحساب أعلى قيمة متكررة حدثت في مجموعة البيانات، ويشير التجميع إلى اتجاه الأسهم من التنبؤ. واستناداً إلى البيانات المالية الرقمية، تم تطوير نموذج جديد للتنبؤ بالأسهم (SPM) لتحليل حركة السوق من خلال مجموعة بيانات رقمية لسوق الأوراق المالية هما S&P 500index و OHLCV. تطوير تقنيات التصنيف المتكاملة التي تجري مع تحليل التنبؤ بناءً على دقة التصنيف حددت بنسبة 82% في هذا البحث وهو أمر واضح انه أفضل من دقة التصنيف في البحوث السابقة. كما تحصل الأداء مع هندسة الميزات في تصنيف النصوص على 93%، في حين أن أداء كل من التصنيف الثنائي والمتعدد المستويات فقد تم دمجها معا للحصول على أفضل مستوى من الدقة. تم تقدير أداء المنهج المقترح من خلال تقييم مختلف المتغيرات كجزء من مجال استرجاع المعلومات، كما تم توضيحه في النتائج التجريبية. ومع ذلك، فإننا لا نوصي باستخدام تطوير تأثيرات النموذج المبني على فلسفة البحوث الأكاديمية من حيث تصنيف البيانات المالية في تحليل التداول الحقيقي.

## **APPROVAL PAGE**

The thesis of Mohammad Rabiul Islam has been approved by the following:

---

Imad Fakhri Al-Shaikhli  
Supervisor

---

Rizal Mohd Nor  
Co-Supervisor

---

Afidalina Tumian  
Co-Supervisor

---

Asadullah Shah  
Internal Examiner

---

Shamala Subramaniam  
External Examiner

---

Hilal Mohammad Yousif Albayatti  
External Examiner

---

Mohamed Elwathig Saeed Mirghani  
Chairman

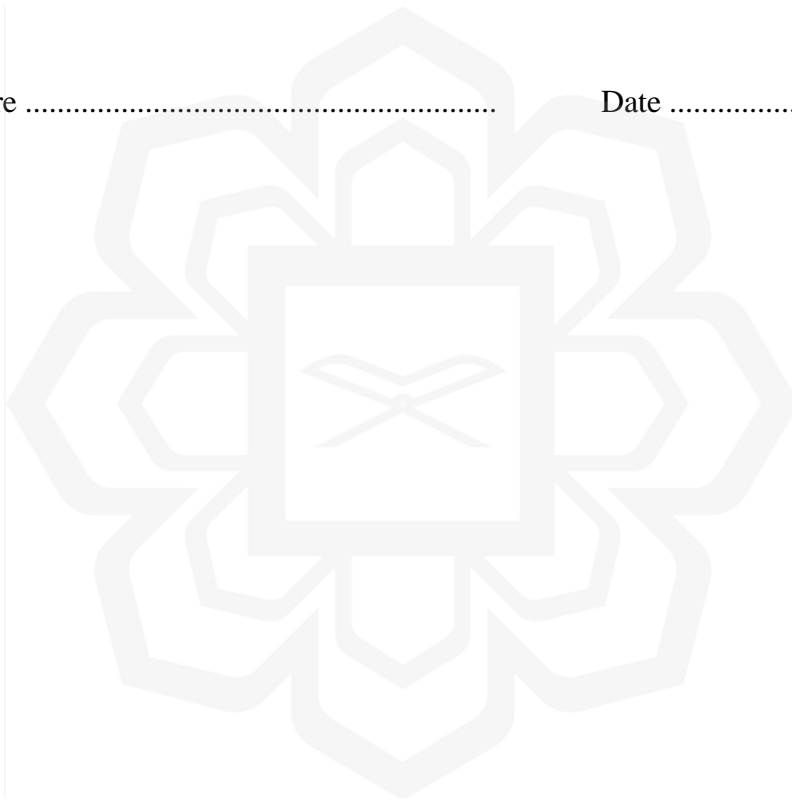
## DECLARATION

I hereby declare that this thesis is the result of my own investigations, except where otherwise stated. I also declare that it has not been previously or concurrently submitted as a whole for any other degrees at IIUM or other institutions.

Mohammad Rabiul Islam

Signature .....

Date .....



**INTERNATIONAL ISLAMIC UNIVERSITY MALAYSIA**

**DECLARATION OF COPYRIGHT AND AFFIRMATION OF  
FAIR USE OF UNPUBLISHED RESEARCH**

**ONLINE FINANCIAL NUMERIC AND TEXTUAL DATA  
ANALYSIS FOR STOCK PREDICTION WITH  
COMPREHENSIVE DATA ANALYSIS THROUGH MACHINE  
LEARNING AND SOFT-COMPUTING TECHNIQUES**

I declare that the copyright holders of this thesis are jointly owned by the student and IIUM.

Copyright © 2020 Mohammad Rabiul Islam and International Islamic University Malaysia. All rights reserved.

No part of this unpublished research may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise without prior written permission of the copyright holder except as provided below

1. Any material contained in or derived from this unpublished research may only be used by others in their writing with due acknowledgement.
2. IIUM or its library will have the right to make and transmit copies (print or electronic) for institutional and academic purposes.
3. The IIUM library will have the right to make, store in a retrieved system and supply copies of this unpublished research if requested by other universities and research libraries.

By signing this form, I acknowledged that I have read and understand the IIUM Intellectual Property Right and Commercialization policy.

Affirmed by Mohammad Rabiul Islam

.....  
Signature

.....  
Date

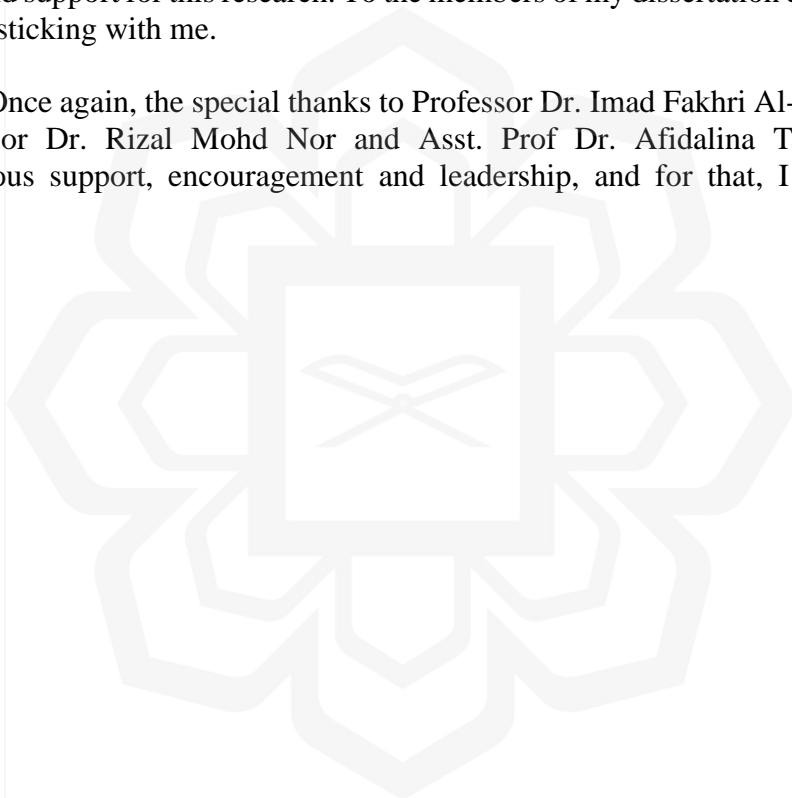
## ACKNOWLEDGEMENTS

Thank you from the bottom of my heart to the Most of all, the Almighty “ALLAH”, the beneficent and the merciful to mankind. Secondly my beloved Professor Dr. Imad Fakhri Taha Al-Shaikhli and rest of the special acquaintance in my life.

Of course, it’s my utmost pleasure to dedicate this work to my dear parents and my family, who granted me the gift of their unwavering belief in my ability to accomplish this goal: thank you for your support and patience.

I wish to express my appreciation and thanks to those who provided their time, effort and support for this research. To the members of my dissertation committee, thank you for sticking with me.

Once again, the special thanks to Professor Dr. Imad Fakhri Al-Shaikhli and co-supervisor Dr. Rizal Mohd Nor and Asst. Prof Dr. Afidalina Tumian for their continuous support, encouragement and leadership, and for that, I will be forever grateful.



# TABLE OF CONTENTS

Abstract .....	ii
Abstract in Arabic .....	iii
Approval Page.....	iv
Declaration .....	v
Copyright .....	vi
Acknowledgements.....	vii
Table of Contents.....	viii
List of Tables .....	xiii
List of Figures .....	xv
List of Abbreviations.....	xviii
<b>CHAPTER ONE: INTRODUCTION .....</b>	<b>1</b>
1.1 Background of the Research .....	1
1.2 Importance of Analysis for Day-Trading Prediction .....	3
1.3 Statement of the Problem.....	3
1.4 Research Questions .....	6
1.5 Research Objectives .....	7
1.6 Conceptual Text Analysis Algorithm Structure.....	8
1.7 Conceptual ANN Model Architecture .....	10
1.8 Research Framework.....	11
1.8.1 Conceptual Framework.....	11
1.8.2 Theoretical Framework.....	13
1.8.2.1 Theoretical Relational Diagram .....	14
1.9 Research Hypotheses .....	17
1.10 Significance of the Research.....	18
1.11 Research Scope .....	20
1.11.1 Portion of Scope.....	20
1.11.2 Research Milestone and Deliverables.....	21
1.12 Research Requirement .....	21
1.13 Thesis Organisation.....	22
1.14 Chapter Summary.....	24
<b>CHAPTER TWO: LITERATURE REVIEW.....</b>	<b>25</b>
2.1 Introduction .....	25
2.2 Literature Review from Theoretical Perspective .....	25
2.2.1 Search Query for this Literature Review .....	27
2.3 Identification of Keywords for Research Gap .....	28
2.4 The Area of Big Data .....	29
2.4.1 Data Mining for Financial Prediction.....	29
2.5 Prediction Techniques of Stock Market.....	30
2.5.1 Importance of Logical Approach for Stock Prediction.....	31
2.5.1.1 Textual and Numerical Data Mining Techniques .....	32
2.5.2 ANN with ARMA Model .....	33
2.5.3 ANN with Back Propagation.....	34
2.5.4 Data Mining with Linear Regression.....	34

2.5.5	Machine Learning Techniques .....	35
2.6	Testing Path for Big Data.....	36
2.7	Mining Techniques and Methods.....	38
2.7.1	Text Mining Techniques.....	38
2.7.1.1	Data Mining Techniques .....	39
2.7.1.2	Opinion Mining Techniques.....	40
2.8	Difference between Text Mining and Numerical Data Mining .....	41
2.8.1	Purpose of Numerical Financial Data Mining .....	42
2.8.2	Purpose of Financial Text Mining .....	43
2.9	Theoretical Review Analysed by Snowball Techniques through Major Relevant Work to Find the Best Solution.....	44
2.10	Chapter Summary.....	50

### **CHAPTER THREE: RESEARCH METHODOLOGY .....51**

3.1	Introduction .....	51
3.2	Proposed Research Methodologies and Selection.....	51
3.3	Selected Main Two Methodologies and Algorithm .....	54
3.3.1	Manipulation for text analysis .....	55
3.3.2	Random Selection for Data Accuracy .....	55
3.3.3	Random Assignment for experiment .....	55
3.3.4	Control for Both Types of Data Analysis .....	56
3.3.5	Text Mining Algorithm Steps .....	56
3.4	Text Processing for Data Dictionary .....	58
3.4.1	Financial Perceptual Computing in Scaling Value.....	58
3.4.2	Score Based Perception Scaling Method.....	59
3.4.3	Human Perception Scaling from Harvard Reference .....	60
3.4.3.1	Preparing the Financial Data Dictionary .....	60
3.4.3.2	Financial Articles Observation by Counting the Positive and Negative Sense of Keywords.....	62
3.4.4	Source of Financial Data Sets.....	65
3.4.5	Text Aggregation from Online and Usable Tools .....	67
3.4.6	Text Pre-Processing .....	69
3.4.7	Data Aggregation Conditional Approach .....	71
3.4.8	Text Mining Return Experiment.....	71
3.5	Textual Data Classification .....	72
3.5.1	Proposed Naïve Bayes Classification .....	72
3.5.2	Naïve Bayes for Binary Classification .....	73
3.5.2.1	Proposed Binary Classification with Data Dictionary.....	74
3.5.2.2	Approach of Binary Classification .....	75
3.5.2.1.1	Binary Searching Techniques.....	76
3.5.2.1.2	Binary Prediction Techniques .....	77
3.5.2.3	Multinomial Classification for Text Analysis .....	78
3.5.2.4	Multinomial Classification with Feature Engineering ....	78
3.5.2.5	Feature Engineering.....	79
3.6	Numerical Data and Aggregation.....	80
3.6.1	Proposed Neural Network With Principal Component Analysis .....	81
3.6.2	Advantages and Drawbacks of PCA.....	82
3.7	The Configuration of Neural Network Architecture.....	83

3.7.1	Weight Matrix of Neural Network .....	86
3.7.2	Operation of the Calculation.....	86
3.7.3	Activation Function for Model Specification .....	86
3.7.4	ReLU Activation Function .....	87
3.7.5	Bias in Neural Network .....	89
3.7.5.1	Cost Function .....	91
3.7.5.2	Neural Network Equation .....	91
3.7.5.3	Important of Scaling Value in PCA.....	93
3.7.5.4	Data Scaling in Dimensionality Reduction .....	93
3.8	Methods of Stock Prediction Model.....	94
3.9	Chapter Summary.....	95

## **CHAPTER FOUR: FINANCIAL TEXT DATA ANALYSIS .....96**

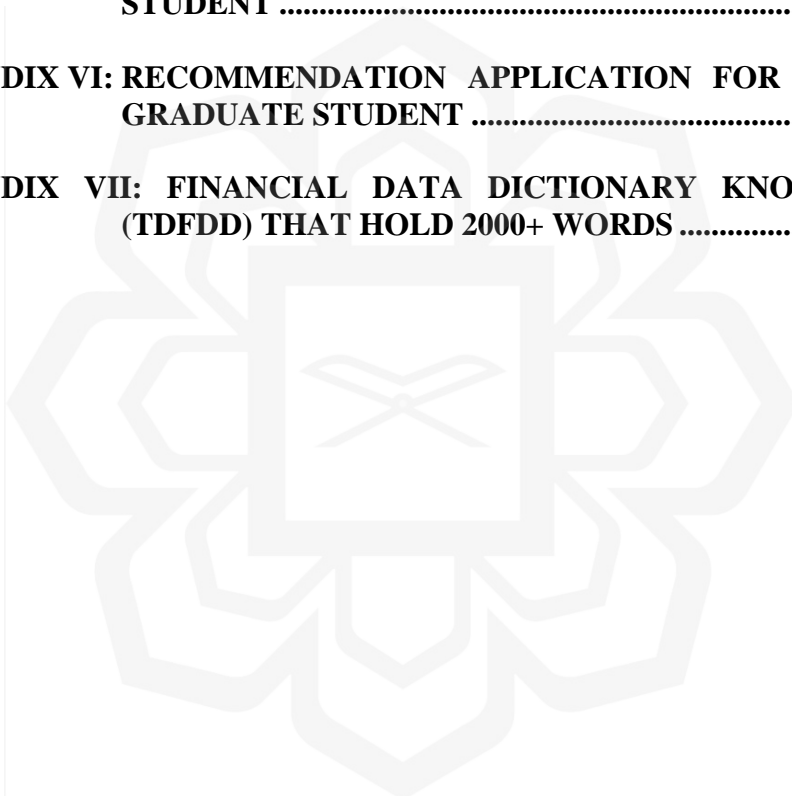
4.1	Introduction .....	96
4.1.1	Interview of Short Survey.....	96
4.1.2	Selected Fuzzy Application Tool .....	99
4.1.3	Hypothesis for Textual Data .....	99
4.2	Identify Population and Sample of Survey .....	100
4.2.1	Survey Objective .....	101
4.2.2	Experiment of Fuzzy Scaling Perception .....	101
4.2.3	Paragraph Scores .....	103
4.2.4	Evaluation of Fuzzy Application on TDFDD.....	106
4.2.4.1	Perception Scaling in Fuzzy Application .....	106
4.2.4.2	Fuzzy Application Test.....	107
4.2.4.3	Hypothesis Proved by Fuzzy Application Test .....	112
4.3	Text Document Analysis Techniques .....	112
4.3.1	Naïve Bayes binary classification on TDFDD .....	113
4.3.2	Text Stemming from Text Documents .....	114
4.3.3	Text Searching Techniques.....	115
4.3.4	Text Analysis and Calculation.....	116
4.4	Financial Text Classification Techniques .....	117
4.4.1	Naïve Bayes Binary Classification .....	118
4.4.2	Binary Classification Operation .....	119
4.4.2.1	Financial Binary Classification Model.....	119
4.4.2.2	Naïve Bayes Binary Classification Model .....	120
4.4.3	Application of Perceptual Computing in Fuzzification .....	125
4.4.3.1	Application of fuzzification in the text parameter.....	125
4.4.4	Binary Classification in IR-Field.....	128
4.5	Naïve Bayes Multinomial Vector Classification.....	131
4.5.1	Word Count Based Dictionary.....	133
4.5.2	Multinomial Classification .....	135
4.6	Chapter Summary.....	137

## **CHAPTER FIVE: FINANCIAL NUMERICAL DATA ANALYSIS.....138**

5.1	Introduction .....	138
5.2	Consequential Stages of Application .....	139
5.3	Structural Dataset .....	143
5.3.1	Database Setup.....	143
5.3.1.1	Preparing Training and Testing Data .....	143

5.3.1.2	Data Scaling with Neural Network.....	144
5.3.1.3	Preparing the Sequence of Setting in Scaling Value.....	145
5.3.1.4	Scaling Value in Principal Component Analysis .....	145
5.3.1.5	Data Scaling with Dimensional Reduction.....	146
5.4	Review of Prediction Techniques .....	147
5.5	ANN Architecture for Stock Data.....	147
5.5.1	Choosing the Best Activation Function.....	149
5.5.2	Neural Network Training Process .....	151
5.5.2.1	Neuron Distribution.....	152
5.5.3	Architecture Model Configuration .....	152
5.5.3.1	Feed Forward-ANN and Calculation Mechanism.....	153
5.5.3.2	Data Inserted to the Neural Network.....	156
5.5.3.3	Train the ANN Model and Performance .....	158
5.5.3.4	Precision and Recall .....	160
5.5.3.5	F-Measure.....	162
5.6.1	Algorithm & Application Field .....	163
5.6.2	Proposed Algorithm VS Existing Algorithm.....	164
5.7	Chapter Summary.....	166
<b>CHAPTER SIX: RESULTS AND DISCUSSION.....</b>		<b>167</b>
6.1	Introduction .....	167
6.2	Experimental Research and Prediction .....	168
6.2.1	Experiment Test and Result .....	169
6.2.2	Pictorial Result Analysis .....	170
6.3	Text Analysis Experiment.....	172
6.3.1	Binary and Multinomial Textual Classification Report.....	172
6.3.2	Binary Classification .....	173
6.3.3	Prediction Through Multinomial Classification .....	176
6.3.4	Numerical Data Classification Analysis and Result.....	179
6.4	Trend Analysis with Top-Down Approach.....	180
6.4.1	Trading Report From Textual Dataset According To Date and Company .....	182
6.4.2	Implication of Numerical Data Analysis and Result for Several Stock Companies .....	184
6.4.3	Implication of Textual Data Analysis and Result for Several Stock Companies.....	186
6.5	Comparison: Previous Result Vs Current Result.....	188
6.6	Chapter Summary.....	189
<b>CHAPTER SEVEN: CONCLUSION AND FUTURE WORK.....</b>		<b>190</b>
7.1	Conclusion .....	190
7.2	Limitation of This Experiment.....	192
7.3	Recommendations and Future Work.....	193

<b>REFERENCES.....</b>	<b>194</b>
<b>APPENDIX I: THE CODE OF ALGORITHM.....</b>	<b>210</b>
<b>APPENDIX II: THE SIMULATION CODE OF THE NAÏVE BAYES ALGORITHM CLASSIFICATION.....</b>	<b>213</b>
<b>APPENDIX III: PUBLICATIONS.....</b>	<b>219</b>
<b>APPENDIX IV: RECOMMENDATION LETTER FOR DATA COLLECTION.....</b>	<b>220</b>
<b>APPENDIX V: LETTER OF APPLICATION FOR POST GRADUATE STUDENT.....</b>	<b>221</b>
<b>APPENDIX VI: RECOMMENDATION APPLICATION FOR UNDER GRADUATE STUDENT.....</b>	<b>222</b>
<b>APPENDIX VII: FINANCIAL DATA DICTIONARY KNOWN AS (TDFDD) THAT HOLD 2000+ WORDS.....</b>	<b>223</b>



## LIST OF TABLES

<u>Table No.</u>		<u>Page No.</u>
1.1	Research Significance in Two Major Sectors	19
1.2	Text Collector from Financial Sources	22
2.1	Identification name and Abbreviation	28
2.2	Comparative Analysis of Text Mining and Data Mining	42
2.3	Cause & Fact for the challenge of Research	45
2.4	Cause & Fact for the Research Gap	45
2.5	Possible Solution for Existing Problem	46
2.6	Current Research and Feasible Solution	47
2.7	Relation Based Possible Solution with the Current Research	47
2.8	Area of Major Research Gap and Problems Based on Various Methods	48
3.1	Article or text-based news sense of score evaluation	63
3.2	Numerical Collection of Stock Market Dataset	65
3.3	Free Sources Based List of Financial News Websites	66
3.4	Paid Sources Based List of Financial News Websites	66
3.5	Advantages and Disadvantages of PCA	83
4.1	Different Data Source collection and Usability	96
4.2	Documents hold the words	135
4.3	Weight calculation of each word	136
5.1	Data Scaling Components	146
5.1.1	Dataset OHLCV	153
5.2	Weight Matrix Given to Four Layers	155
5.3	Confusion Matrix	161
5.4	Algorithm Using on the Different Application Field	163

6.1	Score Components	170
6.2	Accuracy level according to the training set	178
6.3	Intraday Textual Data Analysis for Single Company	183
6.4	Intraday Numerical Data Analysis for Single Company	184
6.5	Precondition for Prediction Result	184
6.6	Intraday Numerical Data Analysis for Multiple Companies	185
6.7	Intraday Textual Data Analysis for Multiple Companies	187
6.8	Previous Models and Current Model Classification Result	188



## LIST OF FIGURES

<u>Figure No.</u>		<u>Page No.</u>
1.1	Consequence Text Mining Steps Extended Trade Analysis	9
1.2	The Sequence of Numerical Datamining Step Based on Proposed Algorithm	10
1.3	Research Variables based on two datasets.	12
1.4	Theoretical Framework Diagram	14
1.5	Text-based term frequency model for financial text data dictionary	16
1.6	The Diagram of Research Significance	19
2.1	Sampling Techniques Diagram for Literature Review	26
2.2	Diagram of Selecting Papers for Literature Review	27
2.3	Text Classification by WordNet	37
3.1	Diagram of Research Methodology	52
3.2	Diagram of Main Two Research Methodologies	54
3.3	Score Conceptual Scale Value from ISV-Scale	59
3.4	Prepared Data Dictionary (TDFDD)	61
3.5	Outwit Hub Pro interface tools	68
3.6	Show the sequence of text pre-processing	69
3.7	Phrase based financial text data	70
3.8	Text Classification Diagram	73
3.9	Flow of TDFDD and News Documents	77
3.10	Excel Sheet of Numerical Daily Basis Stock Data	81
3.11	Neural Network Architectural Diagram	84
3.12	Four Types of Activation Function in Neural Network	87
3.13	ReLU Activation Function	88
3.14	Bias Architecture in Neuron	89

3.15	Calculation between bias and weight values.	90
3.16	Neural Network flow Diagram Based on Equation	92
3.17	Methods of stock prediction model	95
4.1	News Based Survey Questioner for Short Interview Paper	97
4.2	Example of Harvard Online based Dictionary	100
4.3	Diagram shows the perception of the entire experiment	102
4.4	Accumulate total score of positive and negative sense of paragraph	104
4.5	Accumulate the total score in a sense of positive and negative keywords.	105
4.6	Diagram of Intensive Scale Value (ISV)	106
4.7	Depiction of the fuzzy rule-based paragraph in a positive sense of score	108
4.8	Depiction of fuzzy rule-based keywords in positive sense score	109
4.9	Depiction of the fuzzy rule-based paragraph in sense of the negative score	110
4.10	Depiction of fuzzy rule-based keywords in negative sense score	111
4.11	Two types of classification techniques applied to text analysis	112
4.12	TDFDD based Financial Text Analysis	113
4.13	Iteration loop searching techniques based on the condition	115
4.14	Looping Process of Python Code	116
4.15	Condition Based Financial Text Classification	117
4.16	Input Attributes	121
4.17	Financial Document Analysis Scoring Result	128
4.18	Accuracy of the Binary Classification Model	129
4.19	Text Visualization Data Base on Features	130
4.20	Application of Naïve Bayes classification Technique	131
4.21	Count Vectorization from word documents	134

5.1	Excel Sheet of Numerical Daily Basis Stock Data	142
5.2	The flow of Architectural Diagram of Neural Network	148
5.3	Four types of Activation Functions	150
5.4	Architecture Model of NN with Weight Matrix	155
5.5	Input the Numerical Data to Neural Network	157
5.6	Consequence Layer Activities of Neural Network	159
5.7	Conceptualization model of Precision and Recall.	162
5.8	Text Classification Path of Machine Learning Algorithm	164
6.1	S&P 500 index & OHLCV dataset depicts the prediction	171
6.2	Result of text classification analysis	175
6.3	Accuracy result from text dataset	176
6.4	Analysis of Confusion Matrix	178
6.5	Numerical Data analysis result.	179
6.6	Top-Down approach to get the result from two types of datasets	181
6.7	Online based day-trading platform	182
6.8	Comparing the actual result and prediction result	183
6.9	Application of linear growth of explicit equation	187

## LIST OF ABBREVIATIONS

ANN	Artificial Neural Network
ARMA	Auto-Regressive-Moving-Average
BPNN	Back Propagation Neural Network
CNBC	Consumer News and Business Channel
CR	Challenge of Research
CBR	Case Based Reasoning
EMH	Efficient Market Hypothesis
ES	Existing Solution
FN	False Negative
FP	False Positive
FS	Feasible Solution
FFNN	Feed Forward Neural Network
GM	Genetic Algorithm
HTML	Hyper Text Markup Language
ISV	Intensive Scale Value
IR	Information Retrieval
IFR	Financial Information Retrieval
LOCF	Last Observation Carried Forward
LR	Linear Regression
ML-Algo	Machine Learning Algorithm
MLPs	Multilayer Perceptron
MSE	Mean Squared Error
NDTV	New Delhi Television Limited
NMT	News Mining Techniques
NB	Naïve Bayes
NLP	Natural Language Processing
OHLCV	Open, High, Low, Close price and Volume
PTM	Pattern Taxonomy Method
PNN	Probabilistic Neural Networks
PSNR	Pick Signal-to-Noise Ratio
PCA	Principle Component Analysis
PBM	Phrase Based Method
PSP	Potential Solution for Problem
RNN	Recurrent Neural Nets
ReLU	Rectified Linear Unit
RG	Research Gap
SPM	Stock Prediction Model
SVM	Support Vector Machine
S&P 500	Standard and Poor 500 Index
TDFDD	Term Document of Financial Data Dictionary
TBM	Term Based Method
TP	Truth Positive
TN	Truth Negative
WWW	World Wide Web

# CHAPTER ONE

## INTRODUCTION

### 1.1 BACKGROUND OF THE RESEARCH

Stock market prediction is a very active, optimistic and challenging field as to explore in the field of soft computing techniques. The area of financial text mining application is still in its infancy but the potentiality of its future under evolution (Khadjeh Nassirtoussi *et al.*, 2014). A vast amount of big data as collections of the stock market that usually called from the primary source of the World Wide Web (www). In this emerging source of information where a large amount of data able to submerge for the big data system. Having estimated every 20 months, the source of information found double in world wide web (Frawley and Piatetsky-Shapiro, 1992).

Also, such type of information and volume of usable digital data as remained hidden in the web. Data mining and text mining techniques are being applied over the last few years to discover the knowledge from hidden data. (Fayyad, Piatetsky-Shapiro and Smyth, 1996) However, most academicians and practitioners believe that the prediction of the stock market is possible by merely analysing the past stock data. Based on this past numeric or text data, data mining techniques possible to improve by developing some algorithms and applying several strategies to get the prediction in a more effective and efficient way (Raj, 2017).

The technique of text mining enables to extract the high-quality, informative text that usually extracts from a set of textual documents. Analyzing text mining of any field, especially the types of firms that related to information and its application techniques that bring the priority of action-based attributes. These attributes produce the most logical approach instead of low priority (L. Song, 2015). Because the stock market

always deals with a large amount of numeric and non-numeric data that almost impossible to calculate manually in a short time for the traders. In the field of computer science, soft computing techniques utilized the mining techniques as used to predict the stock markets price movement (Oyland, 2015).

Accordingly, advance soft computing techniques such as fuzzy logic in data mining produce many options like operations and logical approaches. Last few decades, various soft computing techniques have been discovered for financial prediction. For example, complex keyword topples that generate the prediction of stock market movements (Wuthrich *et al.*, 1998). Human emotions engaged with hope, fear and worry such critical term come from news that increasing or decreasing amounts of the stock price influence on following day (Zhang, Fuehres and Gloor, 2011). So, this research exactly focus on stock prediction that manifested by soft computing methods.

In this case, “Corpus of Text” was a popular research technique from the last few years that usually contains the word of human emotion, sentiment for data and text mining result in the financial sector (Pang and Lee, 2006; Godbole, Srinivasaiah and Skiena, 2007; Leskovec, Backstrom and Kleinberg, 2009) and many more. Among those techniques, Neural Networks are the most widely used (Choudhry and Garg, 2008). For example, Time Delay Neural Networks for stock market trend prediction, Probabilistic Neural Networks (PNN) used for such model for classification problem. Recurrent Neural Nets (RNN) also used for predicting the next day stock price index with of Bayesian networks, evolutionary algorithms, classifier system, fuzzy algorithm as well as time series analysis (Chung and Shin, 2018) (Jeenanunta, Chaysiri and Thong, 2018) (Camara *et al.*, 2018) (Sharma, 2017) (Environment, 2010).

However, in the financial sector, the significant progress of (Support Vector Machine) SVMs from last several years also impact on financial time series forecasting

(Francis E H Tay & Cao 2001). This beneficial technique impact on technical analysis indicator as input to SVMs. The comparable result of previous BPN and SVM showed that SVM outperformed on BPN more often, but some specific market BPN has also been found better (Sheta, 2006). Instead of BPN and SVM, this research mainly focuses on Feed Forward Neural Network (FFNN) that applies to numerical financial data and Naïve Bayes (NB) for textual day trading data as discussed in subsequent chapters.

## **1.2 IMPORTANCE OF ANALYSIS FOR DAY-TRADING PREDICTION**

Traders of daily basis are known as “day traders” used to active on one-day trading. Due to an intelligent portfolio management system, trading in a day is essential for financial decision making, that directly depend on feasible prediction (Lee, 2007). Moreover, the major stock market indices through the most important financial news that usually taken from daily new closing values as numerical result for decision making (Wuthrich, 1998). Indeed, daily based financial text news brings out the critical decision for day trading. The performance of numeric data or text mining applied on various predicting fields like a policy, engineering, economic, medical, meteorology and others. In this approach, traditional statistical methods also to discovered by hidden knowledge of raw data (Gharehchopogha and Khalifehlou, 2012; Soleimanian, 2012; Soleimanian Gharehchopogh and Reza Khaze, 2012). Throughout this field prediction of algorithm possible to develop in data mining path.

## **1.3 STATEMENT OF THE PROBLEM**

Decision making in investment is very important for the stock market, whether the risk and the problem are in conditional states of a daily market phenomenon due to the vast amount of financial non-numerical data. Previous research found that stock market risk

dependent on daily return observation (Chevapatrakul, Xu and Yao, 2018). The stock price prediction always dependent on technical and fundamental analysis as the expected return comes from different analytical approach (Chaigusin, Chirathamjaree and Clayden, 2008). The technical approach bring on past price, and fundamental financial analysis depends on macroeconomic information, company sector information as well as the company itself (Thomsett, 2005). Both types of information always depend on the open-source sector basis. For example, daily newspaper, online forum, blog, social network, company analytical document that usually published daily, weekly, monthly or even yearly.

Risk is always a common issue as produced from several impacts for lacking the proper situation of the conditional market. A successful trader mostly depends on the contemporary accurate daily based trading price for further trading, but the fluctuation of the stock market is a normal phenomenon. That is why risk-free trading always expected and highly demanding platform for any kinds of financial investment (Prediction, 2016). However, the challenging task is to find the accurate and daily market price from a valuable large amount of textual information, which is usually released online. Most traders depend on daily or weekly price but taking proper decision for price movement gets complicated due to the trade volume in terms of market volatility. Analytical techniques of trade volume able to make a prediction but in some case, data complexity in trade analysis constrained with empirical research that need to extended for simplification (S.R. Dash, 2015) (Chen, 2018).

Nowadays, financial news sources (Numeric and Textual data) are available based on social media like Twitter, Facebook, Yahoo finance, CNN money and many more. Electronic news like textual-data is usually unstructured for testing of any model. This makes weak relationship between online news and one-day stock price movements

due to the noisy data (Noisy data mean other than financial data that cause in this term) experimented in the past (Oyland, 2015). However, this unstructured news is much related to noisy data so these irrelevant noise data need to eliminate from financial text documents that help to get the pure financial text data in terms of financial text-based polar form. The drawback of this research found that the result of accuracy for financial text analysis needs to be improved in terms of classification.

Besides, backpropagation neural network model for numerical data need to improve the architecture by building a new model for stock prediction with feature engineering of the scaling value which integrated principal component analysis (PCA) that enables to improve the accuracy level of stock data classification. The impact of risk defined by the current news and compare it with the actual value of numerical data. Daily basis online news always updated from time to time based on financial news and many researchers active on intelligence portfolio management system, which also recommends daily based stock trading (Lee *et al.*, 2007). The accuracy of prediction also depends on the fluctuation of the market price, which is analyzed by data mining through M-Learning algorithm.

This research points out three problem statements as found in previous betterment of stock-market data analysis.

1. There is a weak relationship between online news and stock volume price movement due to the noisy data and high possibility of variance in succeeding the component of the original data. So, it is inevitable to eliminate noisy data from the original news parameters as input to define more accurate prediction from the financial news (Oyland, 2015; I. Ibidapo, 2017)

2. Online textual news contains the most influential facts or terms for a stock market prediction that hasn't diverse on the dictionary in terms of intraday news-based stock price prediction over the one-day trade volume. This matter does not succeed due to the gap of late arrival news as an input of the text classification model of the same trading day as published (Oyland, 2015; Anbalagan and Maheswari, 2015; P. Gupta, 2017)
3. In the non-linear relationship between the methodology to topology having lack of architectural improvement in terms of weight Matrix value in Artificial Neural Network. That lacking of architectural configuration generate less accurate result with a minimum error of stock data classification (Saini, 2014)

#### **1.4 RESEARCH QUESTIONS**

Research questions bring out the crucial aspect and the analytical thought from unique research questions. As seen in the above research problems, three main questions have pointed out, as stated below, to forward for the research objectives.

1. How do the financial textual news documents overcome the noisy data by eliminating from the original financial textual news?
2. How do the technique and classification model apply to financial textual data to overcome on the intraday news-based stock price prediction over the one-day trade volume that suitable in the gap of late arrival financial news on the same trading day as published?
3. How does architectural neural network perform to minimise the error rate with the efficient algorithm that improves the accuracy level?