

UNIVERSITY OF TECHNOLOGY SYDNEY  
Faculty of Engineering and Information Technology

**MACHINE LEARNING ALGORITHMS FOR  
WEALTH DATA ANALYTICS**

by

**Ngoc Yen Nhi Vo**

A THESIS SUBMITTED  
IN PARTIAL FULFILLMENT OF THE  
REQUIREMENTS FOR THE DEGREE

**Doctor of Philosophy**

Sydney, Australia

2020

## Certificate of Original Authorship

I, Ngoc Yen Nhi Vo, declare that this thesis, submitted in partial fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Computer Science, Faculty of Engineering and IT at the University of Technology Sydney, Australia, is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis. This document has not been submitted for qualifications at any other academic institution. This research is supported by the Australian Government Research Training Program.

Production Note:  
Signature removed  
prior to publication.

Ngoc Yen Nhi Vo  
Sydney, Australia, 2020.

## Acknowledgements

First of all, I express my sincere gratitude for my supervisor, Professor Guandong Xu, who had been a great mentor guiding me through my Ph.D. candidature stages. He had empowered and inspired me to work on multiple research projects which led to publications in top journals and conferences. I am grateful for the support of my co-supervisor, Shaowu Liu, who had been helpful with various technical difficulties. I am also greatly appreciate the contributions of all co-authors to my research publications that partially constructed this thesis: Professor Guandong Xu, Shaowu Liu Ph.D., Professor Xuezhong He, Professor Xitong Li, James Brownlow, Charles Chu, and Ben Culbert.

I thank Advanced Analytics Institute, School of Computer Science, Faculty of Engineering and IT, University of Technology Sydney for providing the infrastructures and computing power used in conducting empirical works. I am deeply grateful for the financial support received during my candidature time: UTS International Research Scholarship, FEIT Industry Scholarship, and Vice Chancellor Conference Fund scholarship. Part of my research was funded by Australian Research Council Linkage Project Scheme under LP170100891 and LP140100937 grants.

Finally, I am grateful for my dear family, especially my father Ngoc Bien Vo, my mother Thi Bao Di Phan, and my brother Ngoc Bao Vo, who had been supportive during my Ph.D. candidature stages, researching and writing my thesis. I am sincerely grateful for everything and especially the opportunities given in this life, from being born to the completion of this Doctor of Philosophy thesis.

# List of Publications

## Journal Papers

- J-1. Vo, Nhi NY, Xuezhong He, Shaowu Liu, and Guandong Xu. "Deep Learning for Decision Making and the Optimization of Socially Responsible Investments and Portfolio." *Decision Support Systems* (2019): 113097.
- J-2. Vo, Nhi NY, Xitong Li, Shaowu Liu, and Guandong Xu. "Leveraging Unstructured Call Log Data for Customer Churn Prediction." *Knowledge-Based Systems* (Under review).

## Conference Papers

- C-1. Vo, Nhi NY, Shaowu Liu, Xuezhong He, and Guandong Xu. "Multimodal mixture density boosting network for personality mining." In *Proceedings of Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 644-655. Springer, Cham, 2018. (PAKDD 2018 conference took place in Melbourne, Australia from 3<sup>rd</sup> to 6<sup>th</sup> June, 2018)
- C-2. Vo, Nhi NY, Shaowu Liu, James Brownlow, Charles Chu, Ben Culbert, and Guandong Xu. "Client Churn Prediction with Call Log Analysis." In *Proceedings of International Conference on Database Systems for Advanced Applications*, pp. 752-763. Springer, Cham, 2018. (DASFAA 2018 conference took place in Queensland, Australia from 21<sup>st</sup> to 24<sup>th</sup> May, 2018)
- C-3. Vo, Nhi NY, and Guandong Xu. "The volatility of Bitcoin returns and its correlation to financial markets." In *Proceedings of International Conference on Behavioral, Economic, Socio-cultural Computing (BESC)*, pp. 1-6. IEEE, 2017. (BESC 2017 conference took place in Cracow, Poland from 16<sup>th</sup> to 18<sup>th</sup> October, 2017)

# Contents

Certificate	ii
Acknowledgments	iii
List of Publications	iv
List of Figures	ix
List of Tables	xii
Abstract	xiv
Abbreviation	xvi
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 Research Objectives . . . . .	3
1.3 Research Highlights . . . . .	4
1.4 Thesis Structure . . . . .	4
<b>2 Theoretical Foundation</b>	<b>6</b>
2.1 Information Retrieval and Data Mining . . . . .	6
2.2 Machine Learning and Algorithms . . . . .	7
2.3 Wealth Data Analytics . . . . .	10
<b>3 Multimodal Mixture Density Boosting Network for Per- sonality Mining</b>	<b>13</b>
3.1 Background and Motivation . . . . .	13

3.2	Preliminary on Personality Mining . . . . .	15
3.3	Multimodal Mixture Density Boosting Network . . . . .	18
3.3.1	DCA Feature Fusion Layer . . . . .	18
3.3.2	Mixture Density Network . . . . .	19
3.3.3	Dynamic Cascade Boosting Network . . . . .	19
3.4	Experiment . . . . .	20
3.4.1	Datasets . . . . .	21
3.4.2	Data Cleaning and Feature Extraction . . . . .	22
3.4.3	Baselines and Evaluation Metrics . . . . .	23
3.5	Result and Discussion . . . . .	24
3.5.1	Empirical Result . . . . .	24
3.5.2	Discussion . . . . .	27
<b>4</b>	<b>Unstructured Data Mining and Interpretable Machine Learning for Wealth Customer Data Analytics</b>	<b>29</b>
4.1	Background and Motivation . . . . .	29
4.2	Preliminary on Wealth Customer Data Analytics . . . . .	33
4.3	Multi-stacking ensemble model for churn prediction . . . . .	37
4.3.1	Unstructured data mining . . . . .	38
4.3.2	Multi-stacking Ensemble Model for Churn Prediction . . . . .	42
4.4	Interpretable Machine Learning for CRM strategies . . . . .	46
4.4.1	SHAP-MRMR+ for interpretable machine learning . . . . .	46
4.4.2	Customer segmentation with Personality . . . . .	47
4.4.3	Customer segmentation with SOM . . . . .	48
4.5	Experiment . . . . .	49

4.5.1	Datasets . . . . .	49
4.5.2	Baselines and Evaluation Metrics . . . . .	54
4.6	Result and Discussion . . . . .	55
4.6.1	Empirical Result . . . . .	55
4.6.2	Robustness Analysis . . . . .	59
4.6.3	Discussion . . . . .	71
<b>5</b>	<b>Data Mining for Socially Responsible Investment</b>	<b>73</b>
5.1	Background and Motivation . . . . .	73
5.2	Preliminary . . . . .	76
5.2.1	Text Mining for Socially Responsible Investment . . . . .	76
5.2.2	ESG as indicator for long-term stock returns forecast . . . . .	77
5.2.3	ESG for Portfolio Optimization . . . . .	78
5.2.4	ESG for Portfolio Diversification . . . . .	79
5.3	Data Mining Methods . . . . .	80
5.3.1	Text Mining Model . . . . .	80
5.3.2	ESG Scores Prediction Model . . . . .	85
5.3.3	P/ESG Indicator Model . . . . .	87
5.3.4	MV-ESG Model . . . . .	88
5.3.5	Combined MV-ESG Model . . . . .	91
5.4	Experiment . . . . .	92
5.4.1	Datasets . . . . .	92
5.4.2	Evaluation Metrics . . . . .	94
5.5	Result and Discussion . . . . .	98
5.5.1	Empirical Result . . . . .	98

5.5.2	Discussion . . . . .	113
<b>6</b>	<b>Deep Learning for Decision Making and Optimization of Socially Responsible Investment Portfolio</b>	<b>116</b>
6.1	Background and Motivation . . . . .	116
6.2	Preliminary . . . . .	119
6.2.1	Socially Responsible Investment . . . . .	119
6.2.2	Deep Learning for Stock Returns Forecasting . . . . .	120
6.2.3	Portfolio Optimization . . . . .	122
6.3	Deep Learning for Socially Responsible Investment Portfolio . . . . .	123
6.3.1	Multivariate BiLSTM for long-term returns prediction . . . . .	123
6.3.2	MV-ESG Multi-Objective Portfolio Optimization . . . . .	127
6.3.3	Reinforcement learning DRIP model . . . . .	129
6.4	Experiment . . . . .	130
6.4.1	Datasets . . . . .	131
6.4.2	Data Cleaning and Feature Extraction . . . . .	131
6.4.3	Baselines and Evaluation Metrics . . . . .	132
6.5	Result and Discussion . . . . .	133
6.5.1	Empirical Result . . . . .	133
6.5.2	Discussion . . . . .	142
<b>7</b>	<b>Conclusion</b>	<b>145</b>
	<b>Bibliography</b>	<b>148</b>
	<b>A Quantifying Socially Responsible Investment Impact</b>	<b>169</b>



## List of Figures

1.1	Shifting Mix of Wealth Data Analytics . . . . .	2
2.1	The Data Mining Process . . . . .	7
2.2	Sample Machine Learning Algorithms . . . . .	9
2.3	Wealth Data Analytics . . . . .	12
3.1	The Five-Factor Model of Personality . . . . .	14
3.2	The complete architecture of our MMDB Neural Network . . . . .	21
3.3	Sample distribution predictions for Openness from MMD network . .	24
4.1	Structured and unstructured data . . . . .	31
4.2	Word Embedding model captures relationships between terms . . . .	40
4.3	Personality Traits Mining Methodology . . . . .	42
4.4	Multi-stacking Ensemble Model and Interpretable Machine Learning .	46
4.5	Method for CRM strategies based on Personalities and Interpretable Machine Learning . . . . .	48
4.6	Example of Kohonen Layer . . . . .	49
4.7	Histograms of LIWC 2015 main features . . . . .	50
4.8	Correlation between Word Embedding Size and AUC Scores . . . . .	55
4.9	Text Features Churn Prediction Model for Investment dataset . . . .	56

4.10 Multi-stacking Ensemble Churn Prediction Model for Investment dataset . . . . .	59
4.11 ROC Curves of prediction models compared with SMOTE methods .	60
4.12 Compare feature impacts with SHAP and SHAP-MRMR+ . . . . .	62
4.13 Compare feature impacts on churn prediction for customer segments with high/low account balance . . . . .	63
4.14 SOM Segments on Employer Superannuation Dataset . . . . .	65
4.15 SOM Segments on Non-Employer Superannuation Dataset . . . . .	66
4.16 SOM Segments with Personalities on Investment Dataset . . . . .	68
4.17 Compare personality impacts on churn prediction for customer A and customer B . . . . .	69
5.1 ESG Rating Framework and Process Overview . . . . .	74
5.3 Our Methodology Framework and Process Overview . . . . .	81
5.4 CSR-Sent Text Feature Extraction Process . . . . .	82
5.5 Standard MV Portfolio with Efficient Frontier . . . . .	89
5.6 Multi-Objective Portfolio Optimization Model . . . . .	91
5.7 Word Cloud of Corporate Social Responsibility Text . . . . .	98
5.8 Plots of stock prices and ESG Scores time series . . . . .	105
5.9 Plots of forecast models evaluation . . . . .	107
5.10 Accumulated Portfolio Values from 2017 to 2018 . . . . .	108
5.11 Pareto Fronts of the “MV-ESG Portfolio” and “MV Portfolio” . . . .	109
5.12 MV and MV-ESG portfolio allocation based on industry . . . . .	111
5.13 MV and MV-ESG portfolio allocation based on negative topics . . . .	112
5.14 MV and Combined MV-ESG portfolio allocation based on industry .	114

6.1	Combined ESG ratings . . . . .	117
6.3	Graphical illustration of LSTM, GRU and BiLSTM . . . . .	126
6.4	Standard MV Portfolio with Efficient Frontier . . . . .	127
6.5	Reinforcement Learning DRIP Model . . . . .	130
6.6	ROC curves . . . . .	136
6.7	MAX-ESG portfolio allocation . . . . .	141
A.1	Investment Impact Measurement Process . . . . .	171

## List of Tables

2.1	Types of Data Analytics . . . . .	11
3.1	Literature Review on Personality Mining . . . . .	17
3.2	10-fold Cross-validation on YouTube and First Impression Datasets .	25
3.3	Transfer learning. . . . .	25
3.4	MMDB and MMD evaluation with YouTube Personality dataset . . .	26
4.1	Existing Text Mining Approaches in Customer Research . . . . .	37
4.2	Sample TFIDF features . . . . .	39
4.3	The Big Five Personality Traits . . . . .	41
4.4	Statistics of the merged DAP datasets . . . . .	52
4.5	Correlation analysis of basic features and customer churn label . . . .	54
4.6	AUC results on the models' prediction accuracy . . . . .	57
4.7	Pair-wise t-test on model performance . . . . .	61
4.8	AUC results with different hyper-parameter settings . . . . .	62
4.9	Predicted churn risk (%) for Customers with Top and Bottom 10% Personality Rank . . . . .	67
5.1	CSR-Sent Text Feature Scores of Intel Corporation . . . . .	83
5.2	Basic statistics of the LIWC-2015 text features . . . . .	84

5.3	Basic statistics of the Empath text features (percentage scores) . . . .	85
5.4	Basic statistics of ESG Ratings dataset . . . . .	93
5.5	Objectives of Tested Portfolios . . . . .	94
5.6	The top correlated text features with Pearson's $r$ . . . . .	100
5.7	10-fold Cross Validation Evaluation of ESG Score Prediction Models .	102
5.8	ESG Scores Correlation Analysis Results . . . . .	104
5.9	Forecast models evaluation . . . . .	106
5.10	P/ESG Scores Correlation Analysis Results . . . . .	108
5.11	MV-ESG Portfolio Evaluation . . . . .	109
5.12	Combined MV-ESG Portfolio Evaluation . . . . .	113
6.1	DRIP Model Evaluation . . . . .	135
6.2	Benchmarking prediction model with multiple hyper-parameters . . .	137
6.3	Benchmarking model with randomly selected datasets . . . . .	138
6.4	MV-ESG Model Evaluation . . . . .	139
6.5	Benchmarking MAX-ESG portfolio with Sustainable Indexes and Funds in 2018 . . . . .	140
6.6	Reinforcement Learning Test Results . . . . .	142

## ABSTRACT

The thesis investigates multiple machine learning algorithms with big data approach and applies cutting-edge deep analytics to tackle the challenges in financial wealth management. In general, the existing research on wealth data analytics is limited with two main challenges. Firstly, the amount of quantitative research conducted is scarce and scattered across different approaches. Partially this is due to the lack of access to the data required for the research to use a quantitative approach. Secondly, the results are rudimentary and limited to a certain aspect of wealth data analytics. This lack of integration in existing research findings is a by-product of the simplistic approaches employed in lieu of big data analytics and deep learning techniques.

This research provides a broader and comprehensive approach for quantitative research within the wealth management field from both financial and customer aspects. Particularly, this research utilizes the big data of structured demographic, behavioral, communicational data, and unstructured textual information from wealth customers, plus additional financial market and corporate responsibility data from companies. This thesis exploits deep analytics techniques to provide a better framework for decision-making support based on the constructed mathematical and computational models, combined with customer segmentation modeling and quantitative finance approach.

From the customer aspect, the thesis applies big data analytics, text mining and interpretable machine learning in customer data analytics in wealth management. The proposed approaches and models are (1) MMDB for personality mining, (2) transfer learning for customer personality prediction, (3) ensemble model with text mining for churn prediction, (4) interpretable machine learning with SHAP-MRMR+ to extract customer insight, and (5) customer segmentation and managerial implications with personality and SOM.

From the financial aspect, this is one of the first research to utilize deep learning

for socially responsible investment. The proposed framework consists of (1) text mining of CSR reports for ESG ratings, (2) ESG-based quantitative models, (3) deep learning using Multivariate BiLSTM for stock return prediction, (4) MV-ESG for ESG-based portfolio optimization, and (5) reinforcement learning for socially responsible investment.

The empirical results show the advantages and effectiveness of deep learning algorithms and big data analytics in financial wealth data analytics. Through the completion of this thesis, various aspects of wealth data analytics have been researched and integrated into sophisticated frameworks, and the information systems can provide meaningful insights for multiple stakeholders, from researchers to individual investors and fund managers.

## Abbreviation

2-D	Two dimensional
AGR	Agreeableness
ARMA	AutoRegressive Moving Average
AUC	Area Under the Curve
BG	Growth of Book Value per Share
Big Five	Big Five Personalities
BiLSTM	Bidirectional Long Short Term Memory network
Borderline-SMOTE	SMOTE and Borderline samples
CL	Call Logs data set
CON	Conscientiousness
CP	Customer Profiles
CRM	Customer Relationship Management
CSR	Corporate Social Responsibility
CSR-Sent	Corporate Social Responsibility Sentiment Dictionary
DCA	Discriminant Correlation Analysis
Doc2Vec	Document to Vector
DRIP	Deep Responsible Investment Portfolio
DT	Decision Tree
DY	Growth of Dividend Yield
EBITDA	Earning Before Interest, Tax, Depreciation and Amortization
EM	Expectation Maximization
ESG	Environmental, Social and Governance
EXT	Extraversion
FI	First Impression Dataset
GARCH	Generalized AutoRegressive Conditional Heteroskedasticity
GHG	Green House Gas



GMO	Genetically Modified Organism
GP	Gaussian Process
GRI	Global Reporting Initiative
GRU	Gated Recurrent Unit network
ID	Identification
IT	Information Technology
K	thousand (quantity)
LAD	Least Absolute Deviation
LI	Lexical Information
LIWC	Linguistic Inquiry and Word Count
LR	Logistic Regression
LSTM	Long Short Term Memory network
MA	Mean Accuracy
MAE	Mean Absolute Error
MAX-ESG	Maximum ESG portfolio
MAX-S	Maximum Sharpe portfolio
MAZ	Mean Absolute Z-Score
MDN	Mixture Density Network
MMD	Multimodal Mixture Density Network
MMDB	Multimodal Mixture Density BoostingNetwork
mRMR	Minimum Redundancy Maximum Relevance
MSE	Mean Squared Error
MV	Mean Variance model
MV-ESG	Mean Variance ESG model
MV-SGP	Mean Variance Stochastic Goal Programming model
NB	Naïve Bayes
NEU	Neuroticism
NLP	Natural Language Processing
NN	Neural Network
NSGA-II	Non-Dominated Sorting Genetic Algorithm 2
OCSVM	One Class Support Vector Machine

OPN	Openness
P/B	Price per Book
P/ESG	Price per ESG
PE	Phrase Embedding
PT	Personality Traits
RF	Random Forest
RL	Reinforcement Learning
RMSE	Root Mean Squared Error
RNN	Recurrent Neural Network
ROC	Receiver Operating Characteristic
SDGs	United Nation Sustainable Development Goals
SHAP	Shapley Addictive Explanations
SHAP-MRMR+	Combined SHAP and positive mRMR method
SLSQP	Sequential Least Square Programming
SMO	Sequential Minimal Optimization
SMOTE	Synthetic Minority Over-sampling Technique
SMOTETomek	SMOTE and Tomek links
SMOTTEENN	SMOTE and Edited Nearest Neighbours
SOM	Self-Organizing Maps
SRI	Socially Responsible Investment
SSS	Small Sample Size
Super	Superannuation
SVD	Single Value Decomposition
SVM	Support Vector Machine
SVM-SMOTE	SMOTE and SVM
SVR	Support Vector Regression
TF-IDF	Term Frequency - Inverse Document Frequency
TI	Term Importance
Word2Vec	Word to Vector
XGB	XgBoost Extreme Gradient Boosting method
YT	Youtube Dataset