

**DECISION MAKING USING
FINANCIAL TIME SERIES DATA**

KARTIKAY GUPTA



**DEPARTMENT OF MATHEMATICS
INDIAN INSTITUTE OF TECHNOLOGY DELHI**

April 2021

© Indian Institute of Technology Delhi (IITD), New Delhi, 2021

**DECISION MAKING USING
FINANCIAL TIME SERIES DATA**

by

KARTIKAY GUPTA

DEPARTMENT OF MATHEMATICS

Submitted

in fulfilment of the requirements of the degree of Doctor of Philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI

April 2021

Dedicated to my family

Certificate

This is to certify that the thesis entitled *Decision Making Using Financial Time Series Data* submitted by *Mr. Kartikay Gupta* to the Indian Institute of Technology Delhi, for the award of the Degree of the Doctor of Philosophy, is a record of the original bona fide research work carried out by him under my supervision and guidance. The thesis has reached the standards fulfilling the requirements of the regulations relating to the degree.

The results contained in this thesis have not been submitted in part or full to any other university or institute for the award of any degree or diploma.



Place: Delhi

Prof. Niladri Chatterjee

Date: 26 April 2021

Supervisor

Department of Mathematics

Indian Institute of Technology Delhi

Acknowledgements

Words are not sufficient to express my gratitude towards people who helped me during the PhD. I had never thought that it would bring so much positivity and openness in my nature.

This PhD would not have been possible without the guidance and encouragement from Prof Niladri Chatterjee, my PhD supervisor. His comments and ideas played a critical role in the PhD. I learned the skill of academic writing through him. He opened up several opportunities in front of me which provided the much-needed exposure.

I remain grateful to my parents Mrs Geeta and Mr Ganesh Parshad for their love and care.

I thank Dr Deepak Joshi for collaborating with me which was a great learning experience. I thank Prof Aparna Mehra for her valuable tips. I am grateful to Prof B.S. Panda, Course coordinator of the course Data Mining and all the students who took this course. I learned many important things regarding PhD research work while working as a TA in this course.

I thank Prof Parag Singla, from whom I learnt significantly by attending two of his courses, namely Machine Learning and Deep Learning. I thank my Student Review Committee faculty members Prof Dharmaraja, Prof Aparna Mehra and Dr Parag Singla for their valuable comments, suggestions and supportive nature.

I also remain filled with gratitude towards all faculty members, PhD scholars and UG students with whom I interacted during my research work. They made my PhD work to be an enjoyable experience.

Place: NIIT Campus, Neemrana

Date: 26 April 2021



Kartikay Gupta

ABSTRACT

Forecasting plays a vital role in the lives of human beings. Time series forecasts are useful for weather forecasting, stocks and derivatives price forecasting, GDP forecasting, electricity load forecasting, among many other things. Good forecasts in any field are immensely valuable. Achieving better forecasts/ predictive-analytics for a stock price-series is the most prominent objective of this work. Specifically, the thesis tries to improvise upon certain different aspects of finance like trading strategies, investment strategies and identifying stock pairs for pairs trading. It works upon techniques of temporal data-mining to propose novel techniques in financial-engineering. Such financial research work helps in extracting money from the market while making the market more efficient and streamlining of the stock prices.

Financial price-series data seldom contains patterns which repeat over time. In this work, first we investigate ways of discovering such patterns (motifs) in financial data and using it for profitable trading. The results indicate the utility of such strategies.

Then, for further improving the forecasts we use fundamental and technical features of companies to predict future prices of their stocks. The technique developed is designed to be scalable to large datasets. The results verify the superiority of the proposed technique over other techniques like simple-regression or random-walk model. The returns generated are much higher than the returns of benchmark indices like NIFTY 50 and BSE Sensex. Here, two large datasets have been used for experimentation.

Our senses can only distinguish between two substances if they are some distance apart in space or time. Measurement of distances, which quantifies farness or nearness between objects, plays an essential role in temporal data-mining or any machine-learning task in general. Thus, the thesis apart from critically reviewing some of the measures, also proposes a few novel distance measures applicable in financial domain.

In this context, we explore clustering for grouping stocks into buckets for quick analysis. Various distance measures have been proposed in the past which can optimally cluster financial time series data. It is generally observed that some stock pairs may not have high Pearson correlation coefficient between them, but they are highly correlated at certain lead-lag. This phenomenon is generally referred to as the lead-lag effect. We propose two new

dissimilarity measures for the above task i.e., Cross-Correlation Type (CCT) and Cross-Correlation Type-II (CCT-II) measure. These measures are able to effectively take into account the time-varying lead-lag relationship between two stocks while measuring their dissimilarity.

We also propose Dynamic Cross-Correlation Type (DCCT) measure which allows the lead-lag relationship between two time series to vary continuously with time. DCCT subtly combines the properties of Dynamic Time Warping (DTW) measure with the CCT measure. Then we showcase the applicability of DCCT measure in identifying pairs for pairs trading. The experiments are conducted to compare the performance of DCCT measure and CCT measure with other popular measures for such task i.e., correlation and SSD measures. The DCCT measure when clubbed with SSD measure i.e., when pairs are selected through optimizing both these measures, then the selected pairs consistently generate the best profit, as compared to all other measures.

Further comparisons show the superiority of the DCCT based lead-lag alignment path over Thermal optimal path (TOP) for determining the lead-lag relationship between the two time series. TOP has been extensively used for analysing lead-lag relationship in the past. This work presents a good alternative for TOP. DCCT measure provides an opportunity to subjectively analyse different price series pairs which have known important lead-lag relationships.

डिसिशन मेकिंग युसिंग फाइनेंसियल टाइम सीरीज डाटा

शोध प्रबंध सार

पूर्वानुमान मनुष्य के जीवन में एक महत्वपूर्ण भूमिका निभाती है। मौसम की भविष्यवाणी, स्टॉक और डेरिवेटिव्स प्राइस फोरकास्टिंग, जीडीपी पूर्वानुमान, बिजली लोड पूर्वानुमान जैसी कई अन्य चीजों के लिए काल श्रेणी का विधिवत अध्ययन उपयोगी होता है। किसी भी क्षेत्र में अच्छे पूर्वानुमान अत्यधिक मूल्यवान होते हैं। शेयर की कीमत का पूर्वानुमान और विश्लेषिकी प्राप्त करना इस कार्य का सबसे प्रमुख उद्देश्य है। विशेष रूप से यह शोध प्रबंध व्यापार के कुछ पहलुओं जैसे निवेश रणनीतियों और जोड़े व्यापार के लिए स्टॉक जोड़े की पहचान करने की रणनीतियों में सुधार करने की कोशिश करता है। इस तरह के वित्तीय शोध कार्य बाजार से पैसे निकालने में मदद करते हैं और साथ ही साथ बाजार को अधिक कुशल बनाते हैं तथा स्टॉक की कीमतों को सुव्यवस्थित करते हैं।

वित्तीय मूल्य-श्रृंखला डेटा में कभी कभी ऐसे पैटर्न (रूपांकनों) होते हैं जो समय के साथ दोहराते हैं। इस काम में, पहले हम वित्तीय आंकड़ों में इस तरह के पैटर्न की खोज करने और लाभदायक व्यापार के लिए इसका उपयोग करने के तरीकों की जांच करते हैं। परिणाम ऐसी रणनीतियों की उपयोगिता का संकेत देते हैं।

फिर हम कंपनियों के मौलिक और तकनीकी विशेषताओं को उनके शेयरों की कीमतों की भविष्यवाणी करने के लिए उपयोग करते हैं। प्रस्तावित तकनीक को बड़े डेटासेट के लिए स्केलेबल बनाया गया है। परिणाम सरल-प्रतिगमन या रैंडम-वॉक मॉडल जैसी अन्य तकनीकों पर प्रस्तावित तकनीक की श्रेष्ठता को सत्यापित करते हैं। उत्पन्न रिटर्न निफ्टी और बीएसई सेंसेक्स जैसे बेंचमार्क सूचकांकों के रिटर्न से अधिक है। यहाँ, प्रयोग के लिए दो बड़े डेटासेट का उपयोग किया गया है।

हमारी इंद्रियां दो पदार्थों के बीच अंतर कर सकती हैं यदि वे अंतरिक्ष या समय में कुछ दूरी पर हैं। दूरियों का मापन, जो वस्तुओं के बीच की मात्रा निर्धारित करता है, सामान्य रूप से डेटा-खनन या किसी मशीन-शिक्षण कार्य में एक आवश्यक भूमिका निभाता है। थिसिस कुछ उपायों की गंभीर रूप से समीक्षा करने के अलावा, वित्तीय क्षेत्र में लागू कुछ दूरी उपायों का भी प्रस्ताव करता है।

इस संदर्भ में, हम विश्लेषण के लिए शेयरों को गुट में समूहित करने के लिए क्लस्टरिंग करते हैं। अतीत में विभिन्न दूरी के उपाय प्रस्तावित किए गए हैं जो वित्तीय समय श्रृंखला डेटा को बेहतर रूप से क्लस्टर कर सकते हैं। यह आमतौर पर देखा गया है कि कुछ स्टॉक जोड़े के बीच उच्च पियरसन सहसंबंध गुणांक नहीं होता है, लेकिन कुछ सीसा-अंतराल पर अत्यधिक सहसंबंध होता है। इस घटना को आम तौर पर लीड-लैग प्रभाव के रूप में जाना जाता है। हम उपरोक्त कार्य के लिए दो नए प्रसार दूरी उपायों का प्रस्ताव करते हैं अर्थात्, क्रॉस-सहसंबंध प्रकार और क्रॉस-सहसंबंध प्रकार-द्वितीय उपाय।

हम डायनेमिक क्रॉस-सहसंबंध प्रकार (डीसीसीटी) उपाय भी प्रस्तावित करते हैं जो दो समय श्रृंखला के बीच लीड-लैग संबंध को समय के साथ लगातार भिन्न करने की अनुमति देता है। फिर हम जोड़े व्यापार के लिए जोड़े की पहचान करने में डीसीसीटी उपाय की प्रयोज्यता का प्रदर्शन करते हैं। इस कार्य के लिए अन्य लोकप्रिय उपायों के साथ तुलना करने के लिए प्रयोग किए गए हैं। डीसीसीटी उपाय जब एक अन्य उपाय के साथ जोड़ा जाता है यानी, जब जोड़े को इन दोनों उपायों को अनुकूलित करने के माध्यम से चुना जाता है, तो चयनित जोड़े लगातार अन्य सभी उपायों की तुलना में सबसे अच्छा लाभ उत्पन्न करते हैं।

आगे की तुलना दो समय श्रृंखला के बीच सीसा-अंतराल संबंध का निर्धारण करने के लिए थर्मल इष्टतम पथ (टीओपी) पर डीसीसीटी आधारित लीड-लैग संरेखण पथ की श्रेष्ठता दिखाती है। बीते समय में सीसा-लैग संबंधों के विश्लेषण के लिए टॉप का बड़े पैमाने पर इस्तेमाल किया गया है। प्रस्तावित काम टीओपी के लिए एक अच्छा विकल्प प्रस्तुत करता है। डीसीसीटी उपाय विभिन्न मूल्य श्रृंखला जोड़े, जिनमें महत्वपूर्ण लीड-लैग संबंध है, का विश्लेषण करने का अवसर प्रदान करता है।

कार्तिकेय गुप्ता

2015MAZ8144

Table of Contents

<i>Certificate</i>	<i>i</i>
<i>Acknowledgements</i>	<i>iii</i>
ABSTRACT	v
<i>List of Tables</i>	<i>xiii</i>
<i>List of Figures</i>	<i>xvi</i>
<i>List of Notations</i>	<i>xviii</i>
Chapter-1	1
1. Introduction	1
1.1 Problem Introduction	2
1.2 Motivation	3
1.3 Thesis Outline	5
Chapter-2	8
2. Multidimensional Motif Discovery in large datasets and forecasting	8
2.1 Motif Discovery Algorithms	8
2.1.1 Univariate time series motif discovery	9
2.1.2 Multivariate time series motif discovery	11
2.1.3 Genetic Algorithms	13
2.2 Proposed Motif-Discovery Algorithm	16
2.2.1 Basic Definitions	16
2.2.2 Description of Genetic Algorithm	20
2.2.3 Genetic Operators	23
2.2.4 Postprocessing of the GA individuals.....	26
2.3 Forecasting through motifs	27
2.4 Experiments	29
2.4.1 Datasets Description	29
2.4.2 Experimental Settings	31
2.4.3 Evaluation Scheme	33
2.4.4 Experiment Results	33

2.5	Conclusion	39
Chapter-3.....		40
3.	<i>Stocks recommendation strategy for higher returns</i>	40
3.1	Introduction	40
3.2	Literature Review.....	41
3.3	Methodology	44
3.3.1	Parameters Training Procedure.....	47
3.3.2	Experimentation Details	48
3.4	Results.....	53
3.5	Discussion.....	59
3.6	Conclusion	60
Chapter-4.....		64
4.	<i>Financial time series clustering into economic sectors</i>	64
4.1	Introduction	64
4.2	Lead-lag relationship.....	65
4.3	Clustering Algorithms.....	66
4.3.1	Agglomerative Hierarchical Clustering	67
4.3.2	K-means clustering.....	69
4.4	Distance Measures.....	70
4.4.1	Correlation Based Dissimilarity Measure (COR).....	72
4.4.2	Temporal Correlation based DTW Dissimilarity Measure (CORT)	72
4.4.3	Shape-Based Distance measure (SBD).....	75
4.5	Proposed Measures	75
4.5.1	Cross-Correlation Type Dissimilarity Measure (CCT)	76
4.5.2	Cross-Correlation Type-II Dissimilarity Measure (CCT-II).....	77
4.6	Experiments Description	78
4.6.1	Methodology.....	79
4.6.2	Datasets description	80
4.6.3	Parameter settings.....	82
4.7	Experiment results.....	82

4.7.1	Indian Data Set	82
4.7.2	S&P500 Data Set	84
4.7.3	DJIA Data Set	84
4.8	Conclusion	86
Chapter-5.....		87
5	<i>Dynamic Lead-Lag relationship measurement between time series</i>	87
5.1	Introduction	88
5.2	Literature Review.....	89
5.3	Distance Measures Description	90
5.3.1	SSD measure.....	90
5.3.2	Pearson Correlation	91
5.3.3	Cross-Correlation Type Measure	91
5.4	Dynamic Cross-Correlation Type measure.....	91
5.5	Methodology and Data description	94
5.5.1	Pairs Trading Strategy	94
5.5.2	Data Description	96
5.6	Results Description	99
5.6.1	Parameter Settings.....	99
5.6.2	Experimental Results	100
5.7	Comparisons with Thermal Optimal Path.....	110
5.8	Discussion.....	117
5.9	Conclusion	119
Chapter-6.....		120
6.	<i>Conclusion</i>	120
6.1	Contributions.....	120
6.2	Future Work	121
Appendix-A.....		123
A.	<i>Feature engineering on temporal data for neuro-degenerative disease classification.....</i>	123

A.1	Feature engineering	123
A.2	Dimensionality reduction & Feature selection	124
A.2.1	Mutual information criteria	125
A.2.2	Random projection	125
A.3	Data Visualisation Tools	125
A.4	Classification Techniques	126
A.4.1	Decision Trees.....	126
A.4.2	Random forests.....	126
A.5	Neuro-degenerative classification Problem	126
A.5.1	Literature Review of the classification problem	127
A.5.2	Data description	129
A.5.3	Feature extraction	132
A.5.4	Classification and Evaluation	135
A.5.5	Results	137
A.5.6	Discussion	140
A.6	Conclusion	145
	<i>Appendix-B</i>	<i>146</i>
B.	<i>Forecasting Algorithms</i>	<i>146</i>
B.1	Random walk model	146
B.2	Vector Autoregression (VAR)	146
B.3	Extreme Gradient boosted decision trees (XGBoost)	147
B.4	SVR	148
	<i>References</i>	<i>151</i>
	<i>About the Author</i>	<i>175</i>

List of Tables

Table 2.1: Range of mutated parameters.	24
Table 2.2: GA Parameters Description.	32
Table 2.3: Overall MAD error rates.	34
Table 2.4: Proposed models results on Thomson Reuters 2D data.	35
Table 2.5: Proposed models results on Thomson Reuters 5D data.	35
Table 2.6: Proposed models results on Synthetic 2D data	36
Table 2.7: Proposed models results on Synthetic 5D data.	37
Table 2.8: XGBoost model results.	38
Table 2.9: SVR model results.	38
Table 3.1: Brief description of different models.	50
Table 3.2: Results Table.	54
Table 3.3: Evaluation scores for BSE data set.	58
Table 3.4: Evaluation scores for NSE data set.	59
Table 3.5: Indicators used for estimating one quarter ahead return for each company.	61
Table 4.1: Hypothetical data set for two time series each of length 37.	77
Table 4.2: The name of the companies whose stock prices are part of the Indian data set.	80
Table 4.3: The 30 Companies comprising the Dow Jones Industrial Average Index.	81
Table 4.4: Cluster evaluation measure for the Indian data set.	84
Table 4.5: Cluster evaluation measure corresponding to the S&P500 data set. ..	84
Table 4.6: Cluster evaluation measure for the DJIA Data Set.	84
Table 5.1: DJIA Company Names.	97
Table 5.2: Sensex 30 Company Names.	98
Table 5.3: TOPIX 30 Company Names.	98
Table 5.4: Parameter settings of DCCT measure (subtable A) and CCT measure (subtable B).	100
Table 5.5: Results corresponding to DJIA dataset with different values of period (P).	102
Table 5.6: Results corresponding to DJIA dataset with different values of period (P).	104

Table 5.7: Results corresponding to Sensex30 dataset with different values of period (P).	106
Table 5.8: Results corresponding to Topix 30 dataset with different values of period (P).	108
Table 5.9: The results of the self-consistency test.	115
Table 5.10: MAD error of the predicted values for different models.	116
Table A.1: Summary of Demographics and severity measures of different groups.	131
Table A.2: The description of various features employed in the decision trees.	139
Table A.3: Sensitivity (Sn), Specificity (Sp) and Accuracy (A) values of the classifier for the classification.	139
Table A.4: Comparisons of the reported approaches.	142

List of Figures

Figure 2.1: A motif in DNA sequences.	9
Figure 2.2: Illustration of a time-series motif.	10
Figure 2.3: A three-dimensional motif in a five dimensional time-series.	13
Figure 2.4: Overview of a basic genetic algorithm.	14
Figure 2.5: A representation of an individual in MOGAMOD.	14
Figure 2.6: A GA individual (b) and its parameter’s description (a).	18
Figure 2.7: A motif structure (b) and its parameter’s description (a).	19
Figure 2.8: The proposed GA.	21
Figure 2.9: Illustration of the mate operator.	22
Figure 2.10: A motif obtained through the proposed GA.	28
Figure 2.11: Overview of the proposed algorithm.	29
Figure 2.12: The 2 motifs inserted into Synthetic 5D data	37
Figure 2.13: A motif in financial time series data.	39
Figure 3.1: Rank-regression Algorithm with regularization.	48
Figure 3.2: AP@100 score for the different models.	55
Figure 3.3: Cumulative Mean Returns.	57
Figure 3.4: Cumulative Mean Returns.	58
Figure 4.1: Brief steps of the hierarchical clustering algorithm.	67
Figure 4.2: A dendrogram formed by hierarchical clustering on financial time series data set.	68
Figure 4.3: Brief steps of the K-means clustering.	69
Figure 4.4: DTW warping path.	74
Figure 4.5: EOD stock prices of the companies	79
Figure 4.6: Dendrogram corresponding to Indian Dataset and dissimilarity measure COR (subfigure a), CORT (subfigure b), CCT (subfigure c) and CCT-II (subfigure d).	83
Figure 4.7: Dendrogram corresponding to DJIA Dataset and dissimilarity measure COR(subfigure a), CORT(subfigure b), CCT (subfigure c) and CCT-II (subfigure d).	85
Figure 5.1: Lead-lag visualization.	93
Figure 5.2: Brief steps of the DCCT measure.	94
Figure 5.3: Brief overview of steps of pairs trading strategy.	95

Figure 5.4: Outline of the overall methodology.....	96
Figure 5.5: Outline of the rolling-window approach.	99
Figure 5.6: Profit Margin by top stock pairs for different measures.....	105
Figure 5.7: Profit Margin by top stock pairs for different measures.....	107
Figure 5.8: Profit Margin by top stock pairs for different measures.....	109
Figure 5.9: Visualisation of the empirically determined lead-lag path.....	113
Figure 5.10: The lead-lag relationship between Verizon-Communications and Coca-Cola companies, determined through DCCT alignment path.	117
Figure A.1: Time series plots for one representative sample from each class. [A - ALS, B - Control, C - Huntington's, D - Parkinson].	132
Figure A.2: Statistical features were extracted from moving windows of different sizes.....	134
Figure A.3: Flow chart representing feature selection methodology as employed in the current chapter.	135
Figure A.4: MI score between features and class labels for different binary classification tasks.....	137
Figure A.5: Decision Tree classifiers obtained for the 4 binary classification tasks.	138
Figure A.6: Confusion matrix of all the classification done.....	140

List of Notations

Symbol	Meaning
\mathbb{R}	Set of real numbers
\mathbb{Z}	Set of integers
\mathbb{C}	Set of complex numbers