

CITY UNIVERSITY OF HONG KONG
香港城市大學

**Data Stream Mining in Financial
Securities Databases**

金融證券數據庫的數據流挖掘

Submitted to
Department of Information Systems
資訊系統學系
in Partial Fulfillment of the Requirements
for the Degree of Doctor of Philosophy
哲學博士學位

by

Liu Xiaoyan
劉曉艷

November 2007
二零零七年十一月

ABSTRACT

In recent years, advances in hardware technology have allowed us to automatically record everyday transactions in stock trading market at a rapid rate. Such processes lead to large amounts of data which grow at an unlimited rate. These data processes are referred to as data streams. One of the key issues for data stream mining is online mining of changes. By understanding the nature of changes in stock data stream, a user may be able to convert this information into valuable trading decisions. Therefore, it is useful to develop tools and techniques which track changes in securities database in real time.

In this research, we focus on the analysis of “T&Q” (Trade and Quote) data, which contains detailed transaction information of securities tick by tick. One data preprocessing technique and several online algorithms to monitor the flow changes in the stock data are presented and their implications in the real securities data are discussed in the thesis.

Since data streams are often in high volume, to reduce the data dimension, we propose a new measure of class heterogeneity and develop a heuristic method to find the approximate optimal discretization scheme in Chapter 3. The numeric evaluation shows that the proposed method can be a good alternative to entropy-based discretization methods.

Segmentation of data streams is useful to find the change in the trend of stock price. In Chapter 4, we propose two novel online segmentation algorithms: the Feasible Window Space method and its extension the Stepwise Feasible Window Space method. They are piecewise-linear-model-based algorithms and always generate fewer segments with acceptable representation error and less computation time. Extensive experiments on a variety of real-world time series are conducted to evaluate the proposed methods.

Monitoring the change in the stock order flow is a meaningful topic in the financial intelligence field since the change in order follow precedes stock price change. In Chapter 5 and Chapter 6, we propose two online change-point detection methods for the stock order flow. One is the multilayer change-point detection algorithm which makes use of the multiresolution property of wavelet transform. It is a non-parametric method. The change-points obtained by this method are more reliable than those detected only from the original time series. The other one is based on the Poisson distribution assumption of the sequence which is a result of empirical finance study. This method identifies the change-points incrementally.

The contributions of this research are two-fold: 1) From the viewpoint of technique, we propose a data reduction method and a set of effective online change detection methods for data streams and develop a series of theoretical results. Detection of changes supports building data mining models in data streams more effectively and accurately. And 2) From the viewpoint of application, the mining of changes in T&Q databases provides a tool to reduce the information asymmetry in the securities market and support the market short-term players making short-term trading decisions.

TABLE OF CONTENTS

ABSTRACT	i
ACKNOWLEDGEMENTS	iii
TABLE OF CONTENTS	iv
LIST OF FIGURES	viii
LIST OF TABLES	ix
CHAPTER 1 INTRODUCTION	1
1.1 Data Stream Mining	1
1.2 Financial Securities Databases.....	3
1.3 Research Questions	5
1.4 Research Objectives and Contributions	8
1.5 Research Methodology	9
1.6 Overview of Dissertation	12
CHAPTER 2 LITERATURE REVIEW	16
2.1 Data Stream Mining	16
2.1.1 The Characteristics of Data Streams.....	16
2.1.2 Literature of Data Stream Mining	18
2.1.2.1 Data Stream Mining Techniques	19
2.1.2.2 Data Stream Management Systems	21
2.1.3 Data Stream Mining in Financial Markets.....	22
2.2 Online Mining of Changes in Data Streams	24
2.2.1 Online Mining of Changes.....	24
2.2.2 Change Mining in Financial Markets	27
2.3 Summary	28

CHAPTER 3 A DISCRETIZATION ALGORITHM BASED ON A HETEROGENEITY CRITERION.....	30
3.1 Introduction to Discretization Techniques.....	30
3.2 Discretization Criterion.....	35
3.2.1 Basic Presentation and Definitions	35
3.2.2 Heterogeneity Discretization Criterion	36
3.3 Heterogeneity-based Discretization Algorithm	45
3.4 Numeric Evaluation	47
3.5 Summary	50
CHAPTER 4 ONLINE SEGMENTATION METHODS FOR DATA STREAMS	52
4.1 Introduction.....	52
4.2 Literature Review on Segmentation.....	54
4.2.1 Segmentation Work	55
4.2.2 Requirements for Segmentation Algorithms	57
4.2.3 Linear Segmentation	58
4.3 Online Segmentation Algorithms	58
4.3.1 Sliding Window Method.....	58
4.3.2 Sliding Window And Bottom-up Method.....	59
4.3.3 Novel Online Segmentation Algorithms: FSW and SFSW	60
4.3.3.1 Maximum Vertical Distance and Segmentation Criterion.....	61
4.3.3.2 Feasible Space Window Method (FSW).....	64
4.3.3.3 Stepwise Feasible Space Window Method (SFSW).....	66
4.4 Complexity Analysis.....	70
4.5 Numeric Evaluation	74
4.5.1 Experiment Setting.....	74
4.5.2 Results on Number of Segments.....	76
4.5.3 Results on Representation Error	77
4.5.4 Overall Performance of Segmentation Methods.....	78

4.5.5 Results on Running Time.....	80
4.5.6 Experiment on Random Walk Series.....	81
4.6 Implication of Segmentation in Stock Data	83
4.6.1 Stock Pattern Design.....	84
4.6.2 Mining Results	85
4.7 Summary	88
CHAPTER 5 MULTILAYER CHANGE-POINT DETECTION IN STOCK ORDER FLOW	90
5.1 Introduction.....	90
5.2 Related Work	93
5.3 Haar Wavelet Transform.....	95
5.3.1 Haar Wavelet Basis.....	95
5.3.2 Discrete Haar Wavelet Transform	97
5.4 Multilayer Change-point Detection Algorithm.....	99
5.4.1 A Meta-algorithm for Change Detection	99
5.4.2 Distance Measure.....	101
5.4.3 Multilayer Change-point Detection Algorithm.....	102
5.5 Numeric Evaluation	105
5.5.1 Poisson Simulation.....	105
5.5.2 Stocks Simulation	111
5.6 Summary	116
CHAPTER 6 INCREMENTAL ONLINE CHANGE-POINT DETECTION IN STOCK ORDER FLOW	118
6.1 Introduction.....	118
6.2 EKOP Sequential Trade Model	119
6.3 Sequential Testing Methods.....	122
6.3.1 Overview of Sequential Tests	122

6.3.2 Sequential Probability Ratio Test (SPRT)	123
6.3.3 Cumulative Sum (CUSUM) Test.....	124
6.4 Incremental Online Change-point Detection Algorithm.....	125
6.4.1 Problem Description	125
6.4.2 Incremental Online Change-point Detection Algorithm.....	126
6.4.3 Algorithm Analysis.....	130
6.4.4 The Extension of the Algorithm	132
6.5 Numeric Evaluation	133
6.5.1 Poisson Simulation.....	134
6.5.2 Stocks Simulation	139
6.6 Summary	140
CHAPTER 7 CONCLUSIONS.....	142
REFERENCES.....	146