

Faculty of Engineering and Information Technology  
University of Technology, Sydney

**Debt Detection and Debt Recovery  
with Advanced Classification  
Techniques**

A thesis submitted in partial fulfillment of  
the requirements for the degree of  
**Doctor of Philosophy**

by

Shanshan Wu

October 2015

## CERTIFICATE OF AUTHORSHIP/ORIGINALITY

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Candidate

---

# Acknowledgments

I would like to express my sincere gratitude to my supervisor, Dr. Ling Chen, for the continuous support of my Ph.D thesis, whose expertise, understanding, and patience added considerably to my thesis completion. I attribute the level of my degree to her encouragement and effort and without her this thesis, too, would not have been completed or written.

My sincere thanks also goes to my co-supervisor, Prof. Chengqi Zhang, for offering me the scholarship opportunities and leading me working on the business project in Centrelink. His advice on both research as well as on my career have been priceless. He has supported me throughout my Ph.D study with his patience and knowledge whilst allowing me the room to work in a flexible time frame.

This thesis was supported by the Australian Research Council (ARC) Linkage Project between University of Technology, Sydney, Australia and Centrelink Australia. I worked on site in Centrelink Canberra for about one and half years. It is one of the most enjoyable working experience I've ever had. A very special thanks goes out to Dr. Yanchang Zhao and Dr. Huaifeng Zhang, who provided me with direction, shared their truthful and illuminating views on a number of issues and became more of a mentor and friend. I must also acknowledge Mr. Hans Bohlscheid, Mr. Peter Newbigin and Mr. Brett Clark from Business Integrity Programs Branch, Centrelink, Australia for their support of domain knowledge and helpful suggestions.

Last but not the least, I would like to thank my family for the support they provided me through my entire life. In particular, I must acknowledge

## *ACKNOWLEDGMENTS*

---

my mum and my husband, without whose love, encouragement and support, I would not have finished this thesis. At the end I would like express appreciation to my beloved son, Jayden, whose birth lights up my life and makes me more powerful to complete the thesis.

# Contents

<b>Certificate</b> . . . . .	<b>i</b>
<b>Acknowledgment</b> . . . . .	<b>ii</b>
<b>List of Figures</b> . . . . .	<b>vii</b>
<b>List of Tables</b> . . . . .	<b>viii</b>
<b>Abstract</b> . . . . .	<b>ix</b>
<b>Chapter 1 Introduction</b> . . . . .	<b>1</b>
1.1 Background . . . . .	1
1.2 Research problem . . . . .	3
1.3 Contributions and structure . . . . .	4
<b>Chapter 2 Related Works</b> . . . . .	<b>7</b>
2.1 Sequence classification . . . . .	7
2.1.1 Feature based classification . . . . .	8
2.1.2 Sequence distance based classification . . . . .	11
2.1.3 Model based classification . . . . .	12
2.2 Fraud Detection . . . . .	13
2.3 Negative Sequential Pattern Mining . . . . .	15
<b>Chapter 3 Adaptive Sequence Classification for Debt Detec-</b> <b>tion</b> . . . . .	<b>17</b>
3.1 Introduction . . . . .	17
3.2 Problem Statement . . . . .	19
3.2.1 Frequent Sequential Patterns . . . . .	19
3.2.2 Classifiable Sequential Patterns . . . . .	19

3.2.3	Discriminative Sequential Patterns . . . . .	20
3.3	Discriminative Frequent Patterns Boosting . . . . .	21
3.4	Adaptive Sequence Classification Framework . . . . .	24
3.5	Case Study . . . . .	25
3.5.1	Data Description . . . . .	26
3.5.2	Effectiveness of Boosting Discriminative Patterns . . . . .	27
3.5.3	Performance of Adaptive Sequence Classification Framework . . . . .	29
3.6	Conclusion and Future Work . . . . .	31
<b>Chapter 4</b>	<b>Debt Detection by Sequence Classification Using Both Positive and Negative Patterns . . . . .</b>	<b>33</b>
4.1	Introduction . . . . .	33
4.2	Sequence Classification Using Both Positive and Negative Sequential Patterns . . . . .	34
4.2.1	Negative Sequential Patterns . . . . .	35
4.2.2	An Algorithm for Building Sequence Classifiers . . . . .	36
4.3	Experimental Results . . . . .	37
4.3.1	Results of Negative Sequential Pattern Mining . . . . .	37
4.3.2	Evaluation of Sequence Classification . . . . .	37
4.4	Conclusions and Discussion . . . . .	39
<b>Chapter 5</b>	<b>Optimising Debt Recovery with Decision Trees . . . . .</b>	<b>45</b>
5.1	Introduction . . . . .	45
5.2	Population and Data Preprocessing . . . . .	47
5.3	A Predictive Model for Customer Ranking . . . . .	48
5.4	Model Evaluation . . . . .	51
5.4.1	Evaluation Criteria . . . . .	51
5.4.2	Comparison with Commercial Software . . . . .	51
5.4.3	Experimental Results on Different Parameters . . . . .	52
5.4.4	Evaluation Results . . . . .	53
5.5	Conclusions and Future Work . . . . .	56

<b>Chapter 6</b>	<b>Conclusions and Future Work</b>	<b>58</b>
6.1	Conclusions	58
6.2	Future work	59
6.2.1	A life-time model for debt detection and prevention	60
6.2.2	Recommendation for debt detection and prevention	62
<b>Chapter A</b>	<b>List of Publications</b>	<b>66</b>
<b>Bibliography</b>		<b>67</b>

# List of Figures

3.1	Architecture of Adaptive Sequence Classification . . . . .	24
3.2	Effectiveness of Discriminative Patterns Boosting . . . . .	28
3.3	RoC curves of Adaptive Sequence Classification Framework . .	31
5.1	Decision Tree . . . . .	49
5.2	Comparison with A Commercial Software (EM) . . . . .	52
5.3	Comparison on Different Parameters . . . . .	53
5.4	Total Voluntary Amount Collected in 12 Weeks (10 runs) . . .	54
5.5	Total Voluntary Amount & Average Voluntary Repayment Amount . . . . .	55
5.6	Total Repayment Amount Collected in 12 weeks (10 runs) . .	56
5.7	Total Repayment Amount & Average Repayment Amount . .	57
6.1	FA of a set of multi-granularity sequences . . . . .	61
6.2	The impacts of various activities . . . . .	63
6.3	Fraud prevention recommendations . . . . .	64



# List of Tables

3.1	Feature-Class Contingency Table . . . . .	21
3.2	Centrelink Data Sample . . . . .	27
3.3	Data Windows . . . . .	30
4.1	Supports, Confidences and Lifts of Four Types of Sequential Rules . . . . .	36
4.2	Selected Positive and Negative Sequential Rules . . . . .	42
4.3	The Pattern Numbers of Each Type in PS10 and PS05 . . . . .	42
4.4	The Classifiers on Various Number of Patterns (PS05-4K) . . . . .	43
4.5	The Classifiers on Various Number of Patterns (PS05-8K) . . . . .	43
4.6	Classification Results with Pattern Set (PS10-4K) . . . . .	44
4.7	The Classifiers on Various Number of Patterns (PS10-8K) . . . . .	44

# Abstract

My study is part of an ARC linkage project between University of Technology, Sydney and Centrelink Australia, which aims to applying data mining techniques to optimise the debt detection and debt recovery. A debt indicates an overpayment made by the government to a customer who is not entitled to that payment.

In social security, an interaction between a customer and the government department is recorded as an activity. Each customer's activities happen sequentially along the time, which can be regarded as a sequence. Based on the experience of debt detection experts, there are usually some patterns in the sequence of activities of customers who commit debts. The patterns indicating the customers' intention to be overpaid can thus be used to discover or predict debt occurrence. The development of debt detection and recovery over sequential transaction data, however, is a challenging problem due to following reasons. (1) The size of transaction data is vast, and the transaction data are being generated continuously as the business goes on. (2) Transaction data are always time stamped by the business system, and the temporal order of the transaction data is highly related to the business logic. (3) The patterns and relationships hidden behind the transaction data may be affected by a lot of factors. They are not only dependent on business domain knowledge, but also subject to seasonal and social factors outside the business. Based on a survey of existing methods on debt detection and recovery, data mining techniques are studied in this thesis to detect and recovery debt in an adaptive and efficient fashion.

Firstly, sequence data is used to model the evolution of customer activities, and the sequential patterns generalize the trends of sequences. For long running sequence classification issues, even if the sequences come from the same source, the sequential patterns may vary from time to time. An adaptive sequential classification model is to be built to make the sequence classification adapt to the sequential pattern variation. The model is applied to 15,931 activity sequences from Centrelink which includes 849,831 activity records. The experimental results show that the proposed adaptive sequence classification framework performs effectively on the continuously arriving data.

Secondly, a new technique of sequence classification using both positive and negative patterns is to be studied, which is able to find the relationship between activity sequences and debt occurrences and also the impact of oncoming activities on the debt occurrence. The same dataset is used for the evaluation. The outcome shows if built with the same number of rules, in terms of recall, the classifier built with both positive and negative rules outperforms traditional classifiers with only positive rules under most conditions.

Finally, decision trees are to be built in the thesis to model debt recovery and predict the response of customers if contacted by phone. The customer contact strategy driven by the model aims to improve the efficiency of debt recovery process. The model is utilized in a real life pilot project for debt recovery in Centrelink. The pilot result outperforms the traditional random customer selection.

In summary, this thesis studies debt detection and debt recovery in social security using data mining techniques. The proposed models are novel and effective, showing potentials in real business.