# Hashing for Large-Scale Structured Data Classification

**Lianhua Chi**

A thesis submitted for the Degree of

Doctor of Philosophy

U|T|S|

Faculty of Engineering and Information Technology

University of Technology, Sydney 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:

Date:

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# ABSTRACT

With the rapid development of the information society and the wide applications of networks, almost incredibly large numbers bytes of data are generated every day from the social networks, business transactions and so on. In such cases, hashing technology, if done successfully, would greatly improve the performance of data management. The goal of this thesis is to develop hashing methods for large-scale structured data classification.

First of all, this work focuses on categorizing and reviewing the current progress on hashing from a data classification perspective.

Secondly, new hashing schemes are proposed by considering different data characteristics and challenges, respectively. Due to the popularity and importance of graph and text data, this research mainly focuses on these two kinds of structured data:

1) The first method is a fast graph stream classification method using Discriminative Clique Hashing (DICH). The main idea is to employ a fast algorithm to decompose a compressed graph into a number of cliques to sequentially extract clique-patterns over the graph stream as features. Two random hashing schemes are employed to compress the original edge set of the graph stream and map the unlimitedly increasing clique-patterns onto a fixed-size feature space, respectively. DICH essentially speeds up the discriminative clique-pattern mining process and solves the unlimited clique-pattern expanding problem in graph stream mining;

2) The second method is an adaptive hashing for real-time graph stream classification (ARC-GS) based on DICH. In order to adapt to the concept drifts of the graph stream, we partition the whole graph stream into consecutive graph chunks. A differential hashing scheme is used to map unlimited increasing features (cliques) onto a fixed-size feature

space. At the final stage, a chunk level weighting mechanism is used to form an ensemble classifier for graph stream classification. Experiments demonstrate that our method significantly outperforms existing methods;

3) The last method is a Recursive Min-wise Hashing (RMH) for text structure. In this method, this study aims to quickly compute similarities between texts while also preserving context information. To take into account semantic hierarchy, this study considers a notion of "multi-level exchangeability", and employs a nested-set to represent a multi-level exchangeable object. To fingerprint nested-sets for fast comparison, Recursive Min-wise Hashing (RMH) algorithm is proposed at the same computational cost of the standard min-wise hashing algorithm. Theoretical study and bound analysis confirm that RMH is a highly-concentrated estimator.