# Complex Graph Stream Mining

Shirui Pan

Faculty of Engineering and Information Technology

University of Technology Sydney

A thesis submitted for the degree of

*Doctor of Philosophy*

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 Student**

_____

# Acknowledgements

I would like to express my earnest thanks to my supervisors, Professor Chengqi Zhang and Professor Xingquan Zhu, who have provided tremendous support and guidance for my research in the past four years. Prof Zhang provided me an opportunity to study in the stimulating and interactive centre for Quantum Computation and Intelligent Systems (QCIS), where I met and leant a lot from many smart and sharp people. I benefit significantly from his unselfish help and invaluable suggestion on my research career. I would like to thank Prof Xingquan Zhu for his continuous guidance and supervisions during my Ph.D. study. Discussing a problem with him has been always a pleasure and eye-opening experience. He always gives me sufficient freedom and encouragement to think and explore my research interest. His vision, creativeness and enthusiasm in solving challenging problems has greatly encouraged me and inspired my works. Without his endless patience, generous support, and constant guidance, this thesis could not have been accomplished.

I would also like to thank all the people that had a positive influence on my day-to-day enjoyment of the job. My office-mates, past and present: Jia Wu, Yifan Fu, Guodong Long, Jing Jiang, Peng Zhang, Tianyi Zhou, Wei Bian, Wei Wang, Xun Wang, Ting Guo, Lianhua Chi, Meng Fang, Mingsong Mao, Zhibin Hong, Hongshu Chen, Shaoli Huang, Haishuai Wang, Mingming Gong, Sujuan Hou, Qin Zhang, Maoying Qiao, Zhiguo Long, Hua Meng, Zhe Xu, Bozhong Liu, Tongliang Liu, Junyu Xuan, and Jiang Bian. They are the ones who have given me support during both joyful and stressful times, to whom I will always be thankful.

Finally, and above all, I want to thank my family for their continuous support. I especially thank my wife, Yu Zheng, who took care of the daily life of our little baby, Yixin Pan, and myself, and shared all my pain, sorrow and joy in every moment of my research. I would like to thank my parents, brothers, and sisters for their unconditional encouragement and support, both emotionally and financially. No words could possibly express my deepest gratitude for their endless love, self-sacrifice and unwavering help. To them I dedicate this dissertation.

# Abstract

Recent years have witnessed a dramatic increase of information due to the ever development of modern technologies. The large scale of information makes data analysis, particularly data mining and knowledge discovery tasks, unprecedentedly challenging. First, data is becoming more and more interconnected. In a variety of domains such as social networks, chemical compounds, and XML documents, data is no longer represented by a flat table with instance-feature format, but exhibits complex structures indicating dependency relationships. Second, data is evolving more and more dynamically. Emerging applications such as social networks continuously generate information over time. Third, the learning tasks in many real-life applications become more and more complicated in that there are various constraints on the number of labelled data, class distributions, misclassification costs, or the number of learning tasks etc.

Considering the above challenges, this research aims to investigate theoretical foundations, study new algorithm designs and system frameworks to enable the mining of complex graph streams from three aspects, including (1) Correlated Graph Stream Mining, (2) Graph Stream Classifications, and (3) Complex Task Graph Classification.

In particular, correlated graph stream mining intends to carry out structured pattern search and support the query of similar graphs from a graph stream. Due to the dynamic changing nature of the streaming data and the inherent complexity of the graph query process, treating graph streams as static datasets is computationally infeasible or ineffective. Therefore, we proposed a novel algorithm, CGStream, to identify correlated graphs from a data stream, by using a sliding window, which covers a number of consecutive batches

of stream data records. Experimental results demonstrate that the proposed algorithm is several times, or even an order of magnitude, more efficient than the straightforward algorithms.

Graph stream classification aims to build effective and efficient classification models for graph streams with continuous growing volumes and dynamic changes. We proposed two methods for complex graph stream classification. Due to the inherent complexity of graph structure, labelling graph data is very expensive. To solve this problem, we proposed a gLSU algorithm, which aims to select discriminative subgraph features with minimum redundancy by using both labelled and unlabelled graphs for graph streams. The second approach handles graph streams with imbalanced class distributions and noise. Both frameworks use an instance weighting scheme to capture the underlying concept drifts of graph streams and achieve significant performance gain on benchmark graph streams.

Complex task graph classification aims to address the graph classification problems with complex constraints. We studied two complex task graph classification problems, cost-sensitive graph classification of large-scale graphs and multi-task graph classification. As in medical diagnosis the misclassification cost/risk for different classes is inherently different and large scale graph classification is highly demanded in real-life applications, we proposed a CogBoost algorithm for cost-sensitive classification of large scale graphs. To overcome the limitation of insufficient labelled graphs for a specific learning task, we further proposed effective algorithms to leverage multiple graph learning tasks to select subgraph features and regularize multiple tasks to achieve better generalization performance for all learning tasks.

# Contents

## II    Graph Stream Classification    65

### Graph Stream Classification: Overview    67

## 5    Graph Stream Classification using Labeled and Unlabeled Graphs    69

## 6    Imbalanced and Noisy Graph Stream Classification    99

# List of Figures

# List of Tables