Faculty of Engineering and Information Technology

University of Technology, Sydney

# Ensemble Predictions: Empirical Studies on Learners' Performance and Sample Distributions

A thesis submitted in fulfillment of

the requirements for the degree of

**Doctor of Philosophy**

by

## Guohua Liang

February 2014

# 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 deepest gratitude to my principal supervisor Professor Chengqi Zhang for providing me with such an excellent research environment; for the wonderful opportunities to explore my dreams to become a successful scientist; for his supervision; ongoing encouragement and support. His deep insights, profound knowledge and outstanding contributions to our society have been a great model for me. All of these will benefit my future research career. It has been a great honor for me to be his PhD student.

I would especially like to express my deepest appreciation to my co-supervisor Professor Xingquan Zhu. His guidance and rigorous approach in all academic areas has been of great value, allowing me to make significant academic achievements in a short period of time. I am very grateful for his support and advice.

My warm thanks to all members, visitors and students of the Centre for Quantum Computation & Intelligent Systems, University of Technology, Sydney for their helpful discussions and encouragement.

I would like to dedicate this thesis to my family, my husband Bin Kong and my two sons, Jeffrey Kong and Kevin Kong for their love, support, and understanding, especially, my little son, Kevin. He always asks me, "Mum why you always go to your office doing your work?" Sometimes he talks to himself, "My mum has a lot of work to do"; he also frequently says to me,

"When you finish your work come back home, I will give you some gold stars."

Finally, I would like to express my gratitude to both the Graduate School and the Faculty of Engineering and IT, University of Technology, Sydney for offering financial support with the final stages of my thesis writing.

# Contents

# List of Figures

# List of Tables

xiv

# Abstract

Imbalanced data problems are among the most challenging in Data Mining and Machine Learning research. This dissertation investigates the performance of ensemble learning systems on different types of data environments, and proposes novel ensemble learning approaches for solving imbalanced data problems. Bagging is one of the most effective ensemble methods for classification tasks. Despite the popularity of bagging in many real-world applications, there is a major drawback on extremely imbalanced data. Much research has addressed the problems of imbalanced data by using over-sampling and/or under-sampling methods to generate an equally balanced training set to improve the performance of the prediction models. However, it is unclear which is the best ratio for training, and under which conditions bagging is outperformed by other sampling schemes on extremely imbalanced data.

Previous research has mainly been concerned with studying unstable learners as the key to ensuring the performance gain of a bagging predictor, with many key factors remaining unclear. Some questions have not been well answered: (1) What are the key factors for bagging predictors to achieve the best predictive performance for applications? and (2) What is the impact of varying the levels of class distribution on bagging predictors on different data environments. There is a lack of empirical investigation of these issues in the literature.

The main contributions of this dissertation are as follows:

1. This dissertation proposes novel approaches, uneven balanced bagging to boost the performance of the prediction model for solving imbalanced problems, and hybrid-sampling to enhance bagging for solving highly imbalanced time series classification problems.

2. This dissertation asserts that robustness and stability are two key factors for building a high performance bagging predictor. This dissertation also derives a new method, utilizing two-dimensional robustness and stability decomposition to rank the base learners into different categories for the purpose of comparing the performance of bagging predictors with respect to different learning algorithms. The experimental results demonstrate that bagging is influenced by the combination of robustness and instability, and indicate that robustness is important for bagging to achieve a highly accurate prediction model.

3. This dissertation investigates the sensitivity of bagging predictors. We demonstrate that bagging MLP and NB are insensitive to different levels of imbalanced class distribution.

4. This dissertation investigates the impact of varying levels of class distribution on bagging predictors with different learning algorithms on a range of data environments, to allow data mining practitioners to choose the best learners and understand what to expect when using bagging predictors.