# Exploring Instance Correlation for Advanced Active Learning

Yifan Fu

A Thesis submitted for the degree of Doctor of Philosophy

Faculty of Engineering and Information Technology

University of Technology, Sydney 2013

# Certificate of Authorship and 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

On having completed this thesis, I am especially thankful to my supervisor Prof. Xingquan Zhu and co-supervisor Prof. Chengqi Zhang, who had led me to an at one time unfamiliar area of academic research, and trusted me and given me as much as possible freedom to purse my own research interests. Prof. Zhu has taught me how to think and study independently and how to solve a difficult scientific problem in flexible but rigorous ways. He has sacrificed much of his precious time for developing my academic research skills. Prof. Zhang has also given me great help and support in life.

I am thankful to the group members I met in the University of Technology, Sydney, including Bin Li, Shirui Pan, Guodong Long, Lianyang Ma, and many others. I learned a lot from these smart people, and I was always inspired by the interesting and in-depth discussions with them. I enjoyed the wonderful atmosphere,being with them, of both academic research and daily life.

I am incredibly grateful to my mother for her generosity and encouragement. This thesis is definitely impossible to be completed without her constant support and understanding. I am also thankful to my friends who have companied me, though not always at my side, through the arduous journey of three and a half years.

# Table of Contents

# List of Figures

# List of Tables

# Abstract

Active learning (AL) aims to construct an accurate classifier with the minimum labeling cost by actively selecting a few number of most informative instances for labeling. AL traditionally relies on some instance-based utility measures to assess individual instances and label the ones with the maximum values for training. However, such approaches cannot produce good labeling subsets. Because instances exist some explicit / implicit relations between each other, instance-based utility measure evaluates instance informativeness independently without considering their interactions. Accordingly, this thesis explores instance correlation in AL and utilizes it to make AL's more accurate and applicable. To be specific, our objective is to explore instance correlation from different views and utilize them for three different tasks, including (1) reduce redundancy for optimal subset selection, (2) reduce labeling cost with a nonexpert labeler and (3) discover class spaces for dynamic data.

First of all, the thesis introduces existing works on active learning from an instance-correlation perspective. Then it summarizes their technical strengths/ weaknesses, followed by runtime and label complexity analysis, discussion about emerging active learning applications and instance-selection challenges therein.

Secondly, we propose three AL paradigms by integrating different instance correlations into three major issues of AL, respectively. 1) The first method is an optimal instance subset selection method (ALOSS), where an expert is employed to provide accurate class labels for the queried data. Due to instance-based utility measures assess individual instances and label the ones with the maximum values, this may result in the redundancy issue in the selected subset. To address this issue, ALOSS simultaneously considers the importance of individual instances and

the disparity between instances for subset selection. 2) The second method introduces pairwise label homogeneity in AL setting, in which a non-expert labeler is only asked "whether a pair of instances belong to the same class ". We explore label homogeneity information by using a non-expert labeler, aiming to further reducing the labeling cost of AL. 3) The last active learning method also utilizes pairwise label homogeneity for active class discovery and exploration in dynamic data, where some new classes may rapidly emerge and evolve, thereby making the labeler incapable of labeling the instances due to limited knowledge. Accordingly, we utilize pairwise label homogeneity information to uncover the hidden class spaces and find new classes timely. Empirical studies show that the proposed methods significantly outperform the state-of-the-art AL methods.