

Multiple Costs and Their Combination in Cost Sensitive Learning

By

Zhenxing Qin

Submitted in fulfilment of the requirement for the degree of

Doctor of Philosophy

University of Technology, Sydney

June 2006

Copyright 2006 by UTS

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

Production Note:

Signature removed prior to publication.....

To my wife Yanfang and Our Parents

Abstract

Cost sensitive learning is firstly defined as a procedure of minimizing the costs of classification errors. It has attracted much attention in the last few years. Being cost sensitive has the strength to handle the unbalance on the misclassification errors in some real world applications. Recently, researchers have considered how to deal with two or more costs in a model, such as involving both of the *misclassification costs* (the cost for misclassification errors) and *attribute test costs* (the cost incurs as obtaining the attribute's value) [Tur95, GGR02, LYWZ04]. Cost sensitive learning involving both attribute test costs and misclassification costs is called *test cost sensitive learning* that is more close to real industry focus, such as medical research and business decision.

Current test cost sensitive learning aims to find an optimal diagnostic policy (simply, a policy) with minimal expected sum of the misclassification cost and test cost that specifies, for example which attribute test is performed in next step based on the outcomes of previous attribute tests, and when the algorithm stops (by choosing to classify). A diagnostic policy takes the form of a decision tree whose nodes specify tests and whose leaves specify classification actions. A challenging issue is the choice of a reasonable one from all possible policies.

This dissertation argues for considering both of the test cost and misclassification cost, or even more costs together, but doubts if the current way, summing up the two costs, is the only right way. Detailed studies are needed to ensure the ways of combination make sense and be “correct”, dimensionally as well as semantically. This

dissertation studies fundamental properties of costs involved and designs new models to combine the costs together.

Some essential properties of attribute test cost are studied. In our learning problem definition, test cost is combined into misclassification cost by choosing and performing proper tests for a better decision. Why do you choose them and how about the ones that are not chosen? Very often, only part of all attribute values are enough for making a decision and rest attributes are left as “unknown”. The values are defined as ‘*absent values*’ as they are left as unknown purposely for some rational reasons when the information obtained is considered as enough, or when patients have no money enough to perform further tests, and so on.. This is the first work to utilize the information hidden in those “absent values” in cost sensitive learning; and the conclusion is very positive, i.e. “Absent data” is useful for decision making. The “absent values” are usually treated as ‘*missing values*’ when left as known for unexpected reasons. This thesis studies the difference between ‘absent’ and ‘missing’. An algorithm based on lazy decision tree is proposed to identify the absent data from missing data, and a novel strategy is proposed to help patch the “real” missing values. .

Two novel test cost sensitive models are designed for different real work scenarios. The first model is a *general test cost sensitive learning framework with multiple cost scales*. Previous works assume that the test cost and the misclassification cost must be defined on the same cost scale, such as the dollar cost incurred in a medical diagnosis. And they aim to minimize the sum of the misclassification cost and the test cost. However, costs may be measured in very different units and we may meet difficulty in defining the multiple costs on the same cost scale. It is not only a technology issue, but

also a social issue. In medical diagnosis, how much money should you assign for a misclassification cost? Sometimes, a misclassification may hurt a patient's life. And from a social point of view, life is invaluable. To tackle this issue, a *target-resource budget learning framework* with multiple costs is proposed. With this framework, we present a test cost sensitive decision tree model with two kinds of cost scales. The task is to minimize one cost scale, called target cost, and keep the other one within specified budgets. To the best of our knowledge, this is the first attempt to study the cost sensitive learning with multiple costs scales.

The second model is based on the assumption that some attributes of an unlabeled example are known before being classified. A test cost sensitive lazy tree model is proposed to utilize the known information to reduce the overall cost. We also modify and apply this model to the batch-test problem: multiple tests are chosen and done in one shot, rather than in a sequential manner in the test-sensitive tree. It is significant in some diagnosis applications that require a decision to be made as soon as possible, such as emergency treatment.

Extensive experiments are conducted for evaluating the proposed approaches, and demonstrate that the work in this dissertation is efficient and useful for many diagnostic tasks involving target cost minimization and resource utilization for obtaining missing information.

Acknowledgements

First, I would like to take this opportunity to express my sincere gratitude to my supervisor, Professor Chengqi Zhang, for his unreserved encouragement, advice and support, and for giving me the opportunity to pursue my PhD at the Faculty of Information Technology at the University of Technology, Sydney. In particular, when I did not meet the enrolment deadline in 2003 due to the IELTS and Visa delay, he made much effort to hold my scholarship for more than half a year until I had the honour of studying and working with him in the past three years, and that is stamped indelibly on my life. His comments are so helpful, and his suggestions are so challenging, that I profit a lot from his guidance during my PhD program. I am honoured and feel happy to be able to work under such a supervisor like Chengqi. His profound knowledge and wisdom have deeply impressed me. I will remember his kindness forever.

Also, I am full of gratitude to Dr. Shichao Zhang, my co-supervisor, for his great help on my daily life and study, and for his detailed and constructive advice on my researches. He is also my master supervisor in China and he recommends me to Professor Chengqi Zhang as his Ph. D student, i.e. he is always a pilot of my life road and guides my onward direction. I would not be able to complete this thesis without his selfless help. He has always provided me with his knowledgeable views.

I am grateful to the Faculty of Information Technology, the University of Technology, Sydney, for providing me with a nice opportunity, an excellent environment and scholarship for the learning and researching here. I am grateful to Dr. Jie Lu and Dr. Guangquan Zhang for their kind support and valuable discussion.

My thanks also go to all those who have helped me in one way or another during my PhD course: Ms. Li Liu, Ms Yanchun Zhou, Mr. Xiaowei Yan, Mr. Longbing Cao, Mr. Qingfeng Cheng, Mr. Chunsheng Li, Mr. Li Lin, Mr. Jiaqi Wang, Mr. Wanli Chen, Mr. Jiarui Ni, Mr. Yanchang Zhao, Mr. Jiarui Ni, Mr. Chengeng Shi, Mr. Xuetao Guo, and Mr. Alan Tao Wang.

Thanks to my wife Yanfang Ji, and our parents. They give me a happy feeling of family.

In addition, I wish to thank the University of California, Department of Information and Computer Science for providing the UCI Repository of machine learning databases [BM98].

List of Publications

The following is a list of my research papers published in the proceedings of referred international conferences or journals during my PhD study at University of Technology, Sydney.

Referred Journal Papers:

- 1 Shichao Zhang, Zhenxing Qin, Charles Ling and Shengli Sheng, "Missing is Useful": Missing Values in Cost-sensitive Decision Trees," **IEEE Transactions on Knowledge and Data Engineering**, Vol. 17 No. 12 (2005): 1689-1693.
- 2 Chengqi Zhang, Zhenxing Qin, Xiaowei Yan, "Association-Based Segmentation for Chinese-Crossed Query Expansion," **IEEE Intelligent Informatics Bulletin** 5(1), 2005: 18-25.

Referred Conference Papers:

- 3 Zhenxing Qin, Chengqi Zhang and Shichao Zhang, "Missing or absent? A Question in Cost-sensitive Decision Tree," **Proceedings of the Fourth International Conference on Active Media Technology (AMT006)**, Jun 2006:
- 4 Yiming Yang, Qiang Yang, Rong Pan et al. and Zhenxing Qin, "Preprocessing Time Series Data for Classification with Application to CRM," In **Proceedings of the 18th Australian Joint Conference on Artificial Intelligence (AI 2005)**, Sydney, Australia, 2005: 133-142.

- 5 Zhenxing Qin, Chengqi Zhang and Shichao Zhang, “Dynamic Test-sensitive Decision Trees with Multiple Cost Scales,” In **Proceedings of International Conference on Fuzzy Systems and Knowledge Discovery (FSKD-2005)**, Changsha, China, August 2005: 402-405.
- 6 Zhenxing Qin, Chengqi Zhang and Shichao Zhang, “Cost-sensitive Decision Trees with Multiple Cost Scales,” In **Proceedings of the 17th Australian Joint Conference on Artificial Intelligence (AI 2004)**, Cairns, Queensland, Australia, 2004: 380-390.
- 7 Zhenxing Qin, Li Liu and Shichao Zhang, “Mining Term Association Rules for Heuristic Query Construction,” In **Proceedings of 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining, (PAKDD 2004)** Sydney, Australia, May 26-28, 2004: 145-154.

Table of Contents

Multiple Costs and Their Combination in Cost Sensitive Learning.....	i
Abstract	iv
Acknowledgements.....	vii
List of Publications.....	ix
Table of Contents	xi
List of Figures.....	xiii
Chapter 1 Introduction	1
1.1 Overview of the Thesis	1
1.2 Contributions	6
1.3 Organization of the Thesis	7
Chapter 2 Background and Literature Review.....	8
2.1 Background and Definition	9
2.1.1 Classic Decision Trees	9
2.1.2 Classic Cost Sensitive Decision Tree.....	12
2.1.3 Types of Cost in Cost Sensitive Learning	17
2.1.4 Unknown Data in Machine Learning and Data Mining.....	21
2.2 Literature Review for Cost-sensitive Learning and Cost Combination.....	23
2.2.1 Classifiers for Single Cost.....	24
2.2.2 Classifiers for Multiple Costs.....	26
2.3 A Cost Sensitive Decision Tree Involving both of Test and Misclassification Cost	29
2.4 Summary	33
Chapter 3 Test Cost Sensitive Decision Trees with Multiple Cost Scales	35
3.1 Motivation	36
3.2 A General Framework for Learning with Multiple Cost Scales	37
3.2.1 Classic cost sensitive learning framework	37
3.2.2 Test cost sensitive learning framework with single cost scale	38
3.2.3 Test cost sensitive learning framework with multiple cost scales	39
3.3 Test Cost Sensitive Learning Decision Tree with Multiple Costs Scales.....	43
3.3.1 Leaf marking criteria.....	44
3.3.2 Attribute selecting criteria for internal nodes	45
3.3.3 Resource Control Issues	48

3.4 Performing Tests on Testing Examples with Resource Control	49
3.5 . Experiments.....	52
3.6 Conclusions and Future Work.....	58
Chapter 4 Utilization based Test Cost Sensitive Decision Trees.....	59
4.1 Lazy Test Cost Sensitive Decision Trees with Multiple Cost Scales	60
4.1.1 Motivation	60
4.1.2 Classic lazy decision tree	60
4.1.3 Lazy Test Cost Sensitive Decision Tree with Two Cost Scales.....	62
4.1.4 Performance Evaluation	65
4.2 Batch Testing Strategies for Test Cost Sensitive Decision Trees	68
4.2.1 Batch Tests Selection	68
4.2.2 Hybrid Lazy Tree for Batch Tests Selection	70
4.2.3 Performance Evaluation	73
4.3 Conclusions and Future Work.....	75
Chapter 5 Absent and Missing Values in Cost-Sensitive Decision Trees	76
5.1 Introduction	77
5.1.1 Missing fields in data set.....	77
5.1.2 Missing or Absent?	78
5.2 Review of Previous Work	81
5.3 Dealing with Missing Values in Cost-sensitive Decision Trees	82
5.3.1 The Known Value Strategy	82
5.3.2 The Null Strategy.....	83
5.3.3 The Internal Node Strategy	84
5.3.4 The C4.5 Strategy	84
5.4 Evaluating and patching up missing values from absent values with hybrid lazy tree	85
5.4.1 Identifying missing data from absent data	86
5.4.2 Patching up missing data.....	86
5.5 Experiments.....	88
5.5.1 Comparing the Four Missing-value Strategies	88
5.5.2 Experiments for identifying absent data	93
5.6 Conclusions and Future Work.....	95
Chapter 6 Conclusions and Future Research.....	97
6.1 Conclusions	97
6.2 Future Research	99
6.2.1 Combination of Multiple costs	100
6.2.2 Properties of other costs	100
Bibliography	102

List of Figures

Figure 2.1	An example of Decision Tree on Credit Card Application.....	10
Figure 2.2	A decision tree built from the Ecoli dataset (costs are set as in Table 2.3).	32
Figure 3.1	Relationship of cost-sensitive learning models	42
Figure 3.2	Three different decision trees for Ecoli data (single cost scale) built with different resource budgets.....	50
Figure 3.3.	Three different decision trees for Ecoli data (multiple cost scales) built with different resource budget	51
Figure 3.4	Comparing the total cost under 3 different resource budgets	52
Figure 3.5.	Comparing the total cost under different resource budgets	55
Figure 3.6.	Comparing the resource utilization (percentage) under different resource budgets	55
Figure 4.1	A generic lazy decision tree algorithm	61
Figure 4.2	Lazy test sensitive decision tree algorithm with two cost scales	64
Figure 4.3.	The total average target costs of single and multiple cost scales tree under different resource budgets (Dataset Ecoli).....	66
Figure 4.4.	The total average target costs for Dataset Breast.	66
Figure 4.5.	The total average target costs for Dataset Heart Disease.....	67
Figure 4.6.	The average total target costs for Dataset Australia.	67
Figure 4.7.	An overall test cost sensitive decision tree	69
Figure 4.8	Lazy test sensitive decision tree algorithm with two cost scales	71
Figure 4.9.	Choosing known attribute as an internal node	72
Figure 4.10	A hybrid lazy decision tree extended from figure 5.2.	72
Figure 4.11.	Total target costs with different ratio of known attributes on dataset Ecoli. Target costs in two strategies go down when more known attributes are available.	73
Figure 4.12.	Total target costs with different ratio of known information on dataset Breast	74

Figure 4.13. Total target costs with different ratio of known information on dataset Heart	74
Figure 4.14. Total target costs with different ratio of known information on dataset Australia	74
Figure 5.1. A decision tree built extended from figure 7.1.....	87
Figure 5.3 Total average costs for Ecoli. In this and the following figures, “KV” stands for the Known Value Strategy, “NULL” for the Null Strategy, “Internal” for the Internal Node Strategy, and “C4.5” for the C4.5 Strategy.....	91
Figure 5.4 Total average costs for Breast	91
Figure 5.5 Total average costs for Heart	92
Figure 5.6 Total average costs for Thyroid.....	92
Figure 5.7 Total average costs for Australia	93
Figure 5.8 Comparing of average cost of three handling strategies on all datasets	94
Figure 5.9 Influence of cost matrix on our patching model	94