

國立交通大學

電機與控制工程學系

碩 士 論 文

2001年資料探勘競賽研究

Study on KDD cup 2001

研 究 生：張文賢

指 導 教 授：林心宇 教授

中華民國九十三年七月

2001年資料探勘競賽研究

Study on KDD cup 2001

研究生：張文賢

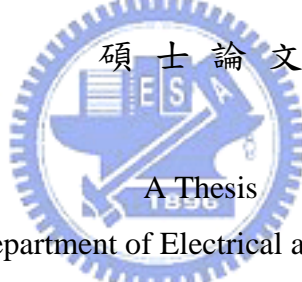
Student: Wen-Hsien Chang

指導教授：林心宇

Advisor: Shin-Yeu Lin

國立交通大學

電機與控制工程學系



Submitted to Department of Electrical and Control Engineering

College of Electrical Engineering and Computer Science

National Chiao-Tung University

in Partial Fulfillment of the Requirements

for the Degree of

Master

in

Electrical and Control Engineering

July 2004

Hsinchu, Taiwan, Republic of China

中華民國九十三年七月


2001年資料探勘競賽研究

學生：張文賢

指導教授：林心宇 博士

國立交通大學電機與控制工程學系（研究所）碩士班

摘 要



資料探勘是一種分析的程序，用來幫助我們發現大型資料庫中的特徵及知識。因為有關生物學的資料探勘快速的發展，2001年資料探勘競賽聚焦在基因及藥物設計資料上。我們所熱衷的是一個分類問題，這個問題有三個有趣的特性（1）大量的遺漏值（2）大量的屬性（3）混合兩種不同型態的資料，而我們最感興趣的分類方法就是決策樹分類法，我們修改了決策樹演算法，並引入“少數服從多數”技巧來提昇分類正確性。為了結合上述兩種分類方法我們發展出“主要-輔助”分類系統。

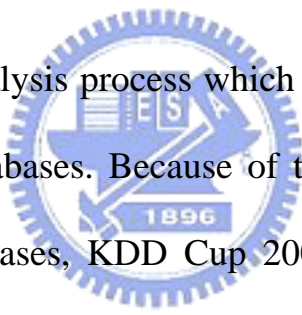
Study on KDD cup 2001

Student: Wen-Hsien Chang

Advisor: Dr. Shin-Yeu Lin

Department of Electrical and Control Engineering
National Chiao Tung University

Abstract

The logo of National Chiao Tung University is a circular emblem. It features a central shield with a book and a torch, surrounded by the university's name in Chinese and English. The year 1896 is inscribed at the bottom of the shield. The entire emblem is encircled by a gear-like border.

Data mining is an analysis process which helps discovering patterns and knowledge in large databases. Because of the rapid growth of interest in mining biological databases, KDD Cup 2001 was focused on data from genomics and drug design. We were involved in a classification problem. The problem has three interesting features: (1) the dataset contains many missing values; (2) this dataset has a lot of attributes; and (3) the dataset is a mixture of two types of data, while the classification method we interested in most is Decision Tree. We modify the Decision Tree algorithm and cite the majority vote to improve the classification accuracy. For integrating the above two classification methods we develop " Primary-Secondary " classification system.

誌 謝

首先感謝指導教授林心宇老師在學習的過程中給予協助與鼓勵，除了讓我在專業領域上有更深刻的體會外，在待人處世和為學態度方面也都獲得相當大的啟發，在此謹致上最誠摯的敬意與謝忱。另外要感謝實驗室士程學長、榮壽學長、佶興學長、志遠學長、傑愷學長提供研究上的意見與經驗，還有學弟紹興在上活上的互相幫忙。最後感謝家人的支持與關懷，讓我能專心順利完成學業，謹以此篇論文獻給你們。



張文賢 于新竹

中華民國九十三年七月

Contents

Abstract (Chinese)	i
Abstract (English)	ii
Acknowledgements	iii
List of tables	vi
List of figures	vii
1 Introduction	1
1.1 The KDD Cup 2001.....	1
1.2 Overview of our study.....	2
2 Problem description and existing approaches	4
2.1 Description of the problem	4
2.1.1 The problem.....	4
2.1.2 The dataset.....	5
2.2 The methods for special features.....	6
2.2.1 The methods for missing values.....	6
2.2.2 The methods for the mass attributes.....	7
2.2.3 The methods for the hybrid type data	8
3 Primary-Secondary classification	9
3.1 Introduction.....	9
3.2 Primary classification.....	10
3.2.1 Primary preprocessing.....	10
3.2.2 Decision tree.....	15
3.2.3 The modified decision tree.....	28
3.3 Secondary classification.....	31
3.3.1 Procedure of secondary classification.....	31

3.3.2	Secondary preprocessing.....	32
3.3.3	Elect the candidate location.....	33
4	What the winners did	34
4.1	Introduction.....	34
4.2	Coping with missing values.....	34
4.3	Nearest neighbor analysis.....	35
4.3.1	Attribute agreement of records.....	35
4.3.2	Neighbors.....	36
4.3.3	Nearest neighbor assignment by prioritizing Attribute.....	36
4.3.4	Classification by Nearest neighbor analysis.....	38
4.4	Computing Optimal Priority.....	38
4.5	Experiment of the winner.....	40
5	Test result and comparison	42
5.1	Primary classification test.....	42
5.1.1	Verification of the classification rules by the modified windowing technique.....	45
5.2	Secondary classification test.....	46
5.3	Overall test result and comparison.....	48
6	Conclusion	50
	Reference	52

List of tables

table 2.1 a transformed training gene.....	4
table 3.1 a training data attribute table.....	12
table 3.2 a transformed training gene G234455.....	15
table 3.3 decision tree keyword	18
table 3.4 pair relation of some genes.....	32
table 3.5 cluster relation of some genes.....	32
table 4.1 interaction relation of gene G234064 andG234126.....	36
table 4.2 Binary interaction relation between pairs of gene.....	36
table 4.3 optimal priority list.....	41
table 5.1 Part of the training data.....	42
table 5.2 classification rules.....	43
table 5.3 example of the ambiguous classification rules.....	46
table 5.4 part of majority voting table (1)	46
table 5.5 part of majority voting table (2)	47
table 5.6 Elect the candidate location by the majority vote	48
table 5.7 The weird attribute values.....	48
table 5.8 comparison of the methods.....	48

List of figures

figure 3.1 Primary-Secondary classification system.....	9
figure 3.2 a decision tree example.....	19
figure 3.3 a decision node.....	20
figure 3.4 output of tree induction algorithm.....	21
figure 3.5 pruning example.....	27
figure 3.6 Nonrecursive decision tree.....	30

