

Faculty of Engineering and Information Technology  
University of Technology, Sydney

# **Coupled Similarity Analysis in Supervised Learning**

A thesis submitted in partial fulfillment of  
the requirements for the degree of  
**Doctor of Philosophy**

by

Chunming Liu

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 Candidate

---

# Acknowledgments

First and foremost, I would like to express my deepest appreciation to my supervisor, Prof. Longbing Cao, for his professional guidance, selfless help and continuous support throughout my PhD study and research. I feel very lucky to have had him as my advisor and I will always remember our many discussions and his invaluable ideas. I sincerely thank him.

I am grateful to my colleagues and friends, Can Wang, Junfu Yin and Xuhui Fan, for their selfless support, and especially to Can Wang, for her advice regarding the technical aspects of my thesis.

My sincere gratitude is extended to my team leader, Zhigang Zheng, for his ongoing support in the Australian Taxation Office project. I am also grateful for the excellent ongoing help that I received from my colleagues and team members, David Wei, Mu Li, and Wei Cao.

I would like to thank all the staff in our Advanced Analytics Institute (AAI). Without their generous support this dissertation would not have been possible.

Last but not the least, I would like to thank my family: my wife and my daughter, for their unconditional love and support throughout my PhD candidature.

Chunming Liu

May 2015 @ UTS

# Contents

Certificate . . . . .	i
Acknowledgment . . . . .	i
List of Figures . . . . .	vi
List of Tables . . . . .	vii
Abstract . . . . .	ix
<b>Chapter 1 Introduction . . . . .</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Limitations and Challenges . . . . .	9
1.3 Research Issues and Objectives . . . . .	14
1.4 Research Contributions . . . . .	15
1.5 Thesis Structure . . . . .	18
<b>Chapter 2 Literature Review . . . . .</b>	<b>21</b>
2.1 Nearest Neighbor Classifier . . . . .	21
2.1.1 $k$ NN . . . . .	22
2.1.2 ROC- $k$ NN . . . . .	24
2.1.3 Fuzzy- $k$ NN . . . . .	25
2.1.4 Summary . . . . .	28
2.2 Similarity for Categorical Data . . . . .	29
2.2.1 Context-free Similarity . . . . .	29
2.2.2 Context-sensitive Similarity . . . . .	34
2.2.3 Summary . . . . .	35
2.3 Class-Imbalance Classification . . . . .	36

## CONTENTS

---

2.3.1	External Methods . . . . .	36
2.3.2	Internal Methods . . . . .	39
2.3.3	Cost-sensitive Methods . . . . .	40
2.3.4	Ensemble Based Methods . . . . .	42
2.3.5	Evaluation . . . . .	44
2.3.6	Summary . . . . .	47
2.4	Multi-Label Classification . . . . .	48
2.4.1	Problem Transformation . . . . .	48
2.4.2	Algorithm Adaptation . . . . .	50
2.4.3	Evaluation . . . . .	52
2.4.4	Summary . . . . .	53
2.5	Summary . . . . .	54
<b>Chapter 3 Coupled <math>k</math>NN for Imbalanced Categorical Data . . . . .</b>		<b>55</b>
3.1	Overview . . . . .	55
3.1.1	Background . . . . .	55
3.1.2	Challenges and Solutions . . . . .	58
3.2	Preliminary Definitions . . . . .	59
3.3	Coupled $k$ NN . . . . .	60
3.3.1	Weights Assignment . . . . .	61
3.3.2	Coupling Similarity . . . . .	62
3.3.3	Integration . . . . .	65
3.3.4	The CF- $k$ NN Algorithm . . . . .	65
3.4	Experiments and Evaluation . . . . .	67
3.4.1	Data and Experimental Settings . . . . .	67
3.4.2	The Performance of CF- $k$ NN . . . . .	68
3.4.3	The Effect of Incorporating Couplings . . . . .	69
3.4.4	The Sensitivity to Imbalance Rate . . . . .	69
3.5	Summary . . . . .	71
<b>Chapter 4 Coupling Based Classification for Numerical Data . . . . .</b>		<b>77</b>
4.1	Overview . . . . .	77

---

4.1.1	Background . . . . .	77
4.1.2	Challenges and Solutions . . . . .	79
4.2	Coupling Relationship on Numerical Attributes . . . . .	80
4.2.1	Problem Statement . . . . .	80
4.2.2	Data Discretization . . . . .	81
4.2.3	Similarity Calculation . . . . .	82
4.2.4	Weight of Coupling . . . . .	83
4.2.5	Integration . . . . .	85
4.3	Experiments and Result . . . . .	86
4.3.1	Data Sets and Settings . . . . .	86
4.3.2	Evaluation Criteria . . . . .	87
4.3.3	Experiments Result . . . . .	88
4.4	Summary . . . . .	90
<b>Chapter 5 Coupled Similarity for Mixed Type Data . . . . .</b>		<b>96</b>
5.1	Overview . . . . .	96
5.1.1	Background . . . . .	96
5.1.2	Challenges and Solutions . . . . .	99
5.2	Preliminary Definitions . . . . .	100
5.3	Coupled Similarity for Mixed Type Data . . . . .	101
5.3.1	Data Discretization . . . . .	101
5.3.2	Weight Calculation . . . . .	102
5.3.3	Similarity Calculation . . . . .	105
5.3.4	Integration . . . . .	107
5.4	Experiments and Evaluation . . . . .	108
5.4.1	Experiments Setting . . . . .	108
5.4.2	Results and Analysis . . . . .	109
5.5	Summary . . . . .	114
<b>Chapter 6 Coupling Analysis in Multi-label Classification . .</b>		<b>121</b>
6.1	Overview . . . . .	121
6.1.1	Background . . . . .	121

*CONTENTS*

---

6.1.2	Challenges and Solutions . . . . .	125
6.2	Methodology . . . . .	126
6.2.1	Problem Statement . . . . .	126
6.2.2	Coupled Label Similarity . . . . .	127
6.2.3	Extended Nearest Neighbors . . . . .	130
6.2.4	Coupled ML- $k$ NN . . . . .	131
6.2.5	Algorithm . . . . .	133
6.3	Experiments and Evaluation . . . . .	134
6.3.1	Experiment Data . . . . .	134
6.3.2	Experiment Setup . . . . .	135
6.3.3	Evaluation Criteria . . . . .	135
6.3.4	Experiment Results . . . . .	136
6.4	Summary . . . . .	137
<b>Chapter 7 Conclusions and Future Work . . . . .</b>		<b>142</b>
7.1	Conclusions . . . . .	142
7.2	Future Work . . . . .	145
<b>Appendix A Appendix: List of Publications . . . . .</b>		<b>147</b>
<b>Appendix B Appendix: List of Symbols . . . . .</b>		<b>149</b>
<b>Bibliography . . . . .</b>		<b>152</b>



# List of Figures

1.1	Categories of Supervised Learning . . . . .	2
1.2	Example of Frequency Vectors . . . . .	6
1.3	Example of Difficulties in Imbalanced Data Sets . . . . .	7
1.4	Examples of Multi-label Images . . . . .	8
1.5	The Profile of Work in This Thesis . . . . .	20
2.1	An Example of 3-Nearest Neighbor Classification . . . . .	23
2.2	Categories of Categorical Data Similarity Measures . . . . .	30
2.3	Example of ROC . . . . .	46
2.4	Categorization of Multi-label Classification Algorithms . . . . .	48
3.1	The Sensitivity of Coupling to Imbalance Rate . . . . .	70
4.1	The Comparison of Specificity. . . . .	88
4.2	The Comparison of Sensitivity. . . . .	89
4.3	The Comparison of Accuracy. . . . .	89
5.1	Sensitivity of IR (CF- $k$ NN: $k$ NN) . . . . .	111
5.2	Sensitivity of IR (CF+ $k$ ENN: $k$ ENN) . . . . .	112
5.3	Sensitivity of IR (CF+CCW $k$ NN:CCW $k$ NN) . . . . .	113

# List of Tables

1.1	Frequency of Co-occurrence . . . . .	9
2.1	Confusion Matrix for A Two-class Problem . . . . .	47
3.1	An Example from The UCI Dataset: Breast Cancer Data . . .	72
3.2	An Example of Frequency of Feature Co-occurrences . . . . .	73
3.3	Data Sets Used in Experiment . . . . .	74
3.4	The AUC Results for CF- $k$ NN in Comparison with Other Algorithms . . . . .	75
3.5	The Comparison of With and Without Coupling . . . . .	76
4.1	Example of Information Table: Wine in UCI . . . . .	92
4.2	Discretization of Information Table: Wine in UCI . . . . .	93
4.3	Numerical Data Sets from UCI . . . . .	94
4.4	The Confusion Matrix of Binary Classification . . . . .	95
5.1	An Fragment from The UCI Dataset: Nursery Data . . . . .	115
5.2	The Frequency of Values Co-occurrence . . . . .	116
5.3	The Data Sets with Mixed Type Features . . . . .	118
5.4	The AUC Results Comparison for HC- $k$ NN and Other Algorithms . . . . .	119
5.5	Comparison for Algorithms With and Without Coupling . . .	120
6.1	An Example of Multi-label Data . . . . .	123
6.2	Transformed Data Sets using Binary Relevance . . . . .	124

6.3	Frequency of Value Pairs . . . . .	128
6.4	CLS Array . . . . .	130
6.5	Extended Nearest Neighbors . . . . .	131
6.6	Experiment Data Sets for Multi-Label Classification . . . . .	139
6.7	Experiment Result1 - Hamming Loss↓ . . . . .	140
6.8	Experiment Result2 - One Error↓ . . . . .	140
6.9	Experiment Result3 - Average Precision↑ . . . . .	141

# Abstract

In supervised learning, the distance or similarity measure is widely used in a lot of classification algorithms. When calculating the categorical data similarity, the strategy used by the traditional classifiers often overlooks the inter-relationship between different data attributes and assumes that they are independent of each other. This can be seen, for example, in the overlap similarity and the frequency based similarity. While for the numerical data, the most used Euclidean distance or Minkowski distance is restricted in each single feature and assumes the features in the dataset have no outer connections. That can cause problems in expressing the real similarity or distance between instances and may give incorrect results if the inter-relationship between attributes is ignored. The same problems exist in other supervised learning, such as the classification tasks of class-imbalance or multi-label. In order to solve these research limitations and challenges, this thesis proposes an insightful analysis on coupled similarity in supervised learning to give an expression of similarity that is more closely related to the real nature of the problem.

Firstly, we propose a coupled fuzzy  $k$ NN to classify imbalanced categorical data which have strong relationships between objects, attributes and classes in Chapter 3. It incorporates the size membership of a class with attribute weight into a coupled similarity measure, which effectively extracts the inter-coupling and intra-coupling relationships from categorical attributes. As it reveals the true inner-relationship between attributes, the similarity strategy we have used can make the instances of each class more compact when

measured by the distance. That brings substantial benefits when dealing with class imbalance data. The experiment results show that our supposed method has a more stable and higher average performance than the classic algorithms.

We also introduce a coupled similar distance for continuous features, by considering the intra-coupled relationship and inter-coupled relationship between the numerical attributes and their corresponding extensions. As detailed in Chapter 4, we calculate the coupling distance between continuous features based on discrete groups. Substantial experiments have verified that our coupled distance outperforms the original distance, and this is also supported by statistical analysis.

When considering the similarity concept, people may only relate to the categorical data, while for the distance concept, people may only take into account the numerical data. Seldom have methods taken into account the both concepts, especially when considering the coupling relationship between features. In Chapter 5, we propose a new method which integrates our coupling concept for mixed type data. In our method, we first do discretization on numerical attributes to transfer such continuous values into separate groups, so as to adopt the inter-coupling distance as we do on categorical features (coupling similarity), then we combine this new coupled distance to the original distance (Euclidean distance), to overcome the shortcoming of the previous algorithms. The experiment results show some improvement when compared to the basic and some variants of  $k$ NN algorithms.

We also extend our coupling concept to multi-label classification tasks. The traditional single-label classifiers are known to be not suitable for multi-label tasks anymore, owing to the overlap concept of the class labels. The most used classifier in multi-label problems, ML- $k$ NN, learns a single classifier for each label independently, so it is actually a binary relevance classifier. As a consequence, this algorithm is often criticized. To overcome this drawback, we introduce a coupled label similarity, which explores the inner relationship between different labels in multi-label classification according

## *ABSTRACT*

---

to their natural co-occurrence. This similarity reflects the distance of the different classes. By integrating this similarity with the multi-label  $k$ NN algorithm, we improve the performance significantly. Evaluated over three commonly used verification criteria for multi-label classifiers, our proposed coupled multi-label classifier outperforms the ML- $k$ NN, BR- $k$ NN and even IBLR. The result indicates that our supposed coupled label similarity is appropriate for multi-label learning problems and can work more effectively compared to other methods.

All the classifiers analyzed in this thesis are based on our coupling similarity (or distance), and applied to different tasks in supervised learning. The performance of these models is examined by widely used verification criteria, such as ROC, Accuracy Rate, Average Precision and Hamming Loss. This thesis provides insightful knowledge for investors to find the inner relationship between features in supervised learning tasks.