# A DIRECT ENSEMBLE CLASSIFIER FOR LEARNING IMBALANCED MULTICLASS DATA

## SAMRY @ MOHD SHAMRIE SAININ

## THESIS SUBMITTED IN FULLFILLMENT FOR THE DEGREE OF DOCTOR OF PHILOSOPHY

## SCHOOL OF ENGINEERING AND INFORMATION TECHNOLOGY
## UNIVERSITI MALAYSIA SABAH
## 2013

# UNIVERSITI MALAYSIA SABAH

BORANG PENGESAHAN STATUS TESIS

JUDUL :     A DIRECT ENSEMBLE CLASSIFIER FOR LEARNING IMBALANCED
            MULTICLASS DATA

IJAZAH:     DOKTOR FALSAFAH (SAINS KOMPUTER)

Saya  <u>SAMRY @ MOHD SHAMRIE SAININ,</u>  Sesi pengajian 2009-2013, mengaku
membenarkan tesis ini disimpan di Perpustakaan Universiti Malaysia Sabah dengan
syarat-syarat kegunaan seperti berikut:-

1.  Tesis ini adalah hak milik Universiti Malaysia Sabah.
2.  Perpustakaan Universiti Malaysia Sabah membenarkan membuat salinan
    untuk tujuan pengajian sahaja.
3.  Perpustakaan dibenarkan untuk membuat salinan tesis ini sebagai bahan
    pertukaran antara institusi pengajian tinggi.
4.  Sila tandakan ( / )

<br>

☐  SULIT        (Mengandungi maklumat yang berdarjah keselamatan
                atau kepentingan Malaysia seperti yang termaktub di
                dalam AKTA RAHSIA RASMI 1972)

☐  TERHAD       (Mengandungi maklumat TERHAD yang telah
                ditentukan oleh organisasi/badan di mana penyelidikan
                dijalankan)

☑  TIDAK TERHAD

<br>

Disahkan oleh

<br>

_____                    _____
(Tandatangan Penulis)                       (Tandatangan Pustakawan)
                                            NURULAIN BINTI ISMAIL
                                            LIBRARIAN
                                            UNIVERSITI MALAYSIA SABAH

<br>

                                            _____
                                            (DR. RAYNER ALFRED)
Tarikh: 20 Disember 2013                     Penyelia

# DECLARATION

I hereby declare that the material in this thesis is my own except for quotations, excerpts, equations, summaries and references, which have been dully acknowledged.

6 September 2013

Samry @ Mohd Shamrie Sainin
PK20088357

# CERTIFICATION

NAME            : **SAMRY @ MOHD SHAMRIE SAININ**

MATRIC NO.      : **PK20088735**

TITLE           : **A DIRECT ENSEMBLE CLASSIFIER FOR LEARNING IMBALANCED MULTICLASS DATA**

DEGREE          : **DOCTOR OF PHILOSOPHY (COMPUTER SCIENCE)**

VIVA DATE       : **17 MAY 2013**

## DECLARED BY

**1. SUPERVISOR**
   Dr. Rayner Alfred

Signature

iii

# ACKNOWLEDGEMENT

First and foremost, in the name of Allah the Most Gracious and the Most Merciful. I would like to extend my deepest praise to Allah S.W.T for the blessings and guidance in completing this research and thesis.

I would like to express my deepest appreciation to my supervisor, Dr. Rayner Alfred for all his advices, guidance and support in this interesting research work that lead to the completion of this thesis. You are the best supervisor.
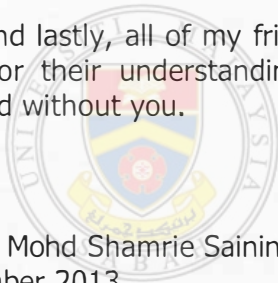
The Ministry of Higher Learning and Universiti Utara Malaysia for the financial support throughout my studies, and also Universiti Malaysia Sabah for the resources and facilities provided.

My wife, Suraya Alias, my sons, Aiman Haziq, new born Aish Hafiy and daughter Damia Hana, for the love, trust, inspiration and great understanding that had given me.

My parents Sainin Ghane and Rosnah Moringking, and my family for being supportive towards my studies.

And lastly, all of my friends whether in the academic circle or outside of my studies for their understanding, advices and support. This research will not be completed without you.

Samry @ Mohd Shamrie Sainin
6 September 2013

# ABSTRACT

A traditional direct single classifier can be easily applied to solve a multiclass classification problem. However, the performance of a single classifier is decreased with the existence of imbalanced data in multiclass classification tasks. Thus, an ensemble of classifiers is one of the methods used to solve multiclass classification tasks. In this thesis, the problem of learning from imbalanced multiclass data classification is studied. In the multiclass classification problem, decision can be estimated not only by the final single class label, but also by other appropriate class. Many real-world multiclass classification problems can be represented into a setting where non-crisp label need to be observed. An in-depth review and method to solve this special learning task is explained in this thesis. An alternative ensemble learning framework called Direct Ensemble Classifier for Imbalance Learning (DECIML) is proposed combining the advantages of existing single classifiers and ensemble methods and strategies. The learning framework consists of ensemble learning and decision combiner model with general supervised learning algorithms as base learner. Feature selection is also applied in DECIML in order to increase the performance of the ensemble learning. In order to facilitate the experiments and future research on the imbalanced multiclass problem, a standard pool of benchmark data is created, which consists of 16 datasets with different degrees of imbalanced ratio and 4 datasets for imbalanced multiclass with feature selection purposes. The benchmark data is used to evaluate and compare the proposed frameworks with several ensemble methods, such as bagging and adaboost. The DECIML with feature selection is also evaluated and compared with methods named CFsSubsetEval and Filteredsubseteval. The results obtained show that the proposed learning frameworks are comparable to other methods. In addition, the selected benchmark data, experiments and the results are useful for future research on the imbalanced multiclass classification problem. Furthermore, the DECIML framework was applied to the real world leaf classification problem based on the shape features. Extensive experiments and results show that the DECIML method does provide a promising performance in imbalanced multiclass with highly noisy data.

# ABSTRAK

## A DIRECT ENSEMBLE CLASSIFIER FOR LEARNING IMBALANCED MULTICLASS DATA

*Algoritma pengelasan tunggal tradisional boleh digunakan dengan mudah secara langsung untuk pelbagai masalah klasifikasi berbilang-kelas. Walau bagaimanapun, prestasi pengelas tunggal akan menurun dengan kewujudan ketidakseimbangan dalam tugas klasifikasi berbilang kelas. Oleh itu, kombinasi pengelas adalah salah satu kaedah dalam tugas klasifikasi berbilang-kelas untuk masalah ketidakseimbangan dalam perlombongan data dan pembelajaran mesin. Dalam tesis ini, masalah pembelajaran dari klasifikasi data berbilang-kelas tidak seimbang dikaji. Dalam masalah pengelasan berbilang-kelas, keputusan boleh dianggarkan bukan sahaja oleh label kelas akhir tunggal, tetapi kelas yang sesuai yang lain. Klasifikasi masalah berbilang-kelas dalam dunia sebenar kebanyakannya boleh diwakilkan menggunakan label bukan tunggal yang perlu dipatuhi. Suatu kajian semula yang mendalam dan kaedah untuk menyelesaikan tugas pembelajaran khas dijelaskan dalam disertasi ini. Rangka kerja alternatif bagi kombinasi pembelajaran yang dikenali sebagai Kombinasi Pengelas Pembelajaran Ketidakseimbangan Berbilang Kelas Secara Langsung (DECIML) dicadangkan berdasarkan kepada kelebihan pengelas tunggal yang sedia ada dan kaedah serta strategi kombinasi. Rangka kerja pembelajaran ini terdiri daripada kombinasi pembelajaran dan penggabung keputusan model dengan algoritma pembelajaran terselia sebagai pembelajar asas. Satu lagi rangka kerja pembelajaran ialah menggabungkan DECIML dan pemilihan ciri untuk meningkatkan prestasi kombinasi pembelajaran. Bagi memudahkan ujikaji dan kajian akan datang untuk data berbilang-kelas tidak seimbang, satu senarai piawai data sebagai tanda aras diwujudkan, yang mana terdiri daripada 16 set data dengan darjah nisbah ketidakseimbangan yang berbeza dan 4 dataset untuk berbilang-kelas tidak seimbang untuk tujuan pemilihan ciri. Data penanda aras ini digunakan untuk menilai dan membandingkan rangka kerja yang dicadangkan dengan beberapa kaedah kombinasi, seperti bagging dan adaboost. DECIML dengan pemilihan ciri juga dinilai dan dibandingkan dengan kaedah seperti CFsSubsetEval dan Filteredsubseteval. Hasil kajian menunjukkan bahawa rangka kerja pembelajaran yang dicadangkan adalah setanding dengan kaedah lain. Di samping itu, data tanda aras yang dipilih, eksperimen dan keputusan boleh digunakan untuk penyelidikan masa depan dalam masalah klasifikasi berbilang-kelas tidak seimbang. Di samping itu, rangka kerja DECIML telah digunakan untuk klasifikasi dunia sebenar, masalah klasifikasi daun berdasarkan ciri-ciri bentuk. Ujikaji yang mendalam dan keputusan yang diperolehi menunjukkan bahawa kaedah DECIML memberikan prestasi yang baik dalam masalah berbilang-kelas tidak seimbang dengan data yang sangat bising. Oleh itu, penemuan menarik daripada keputusan eksperimen adalah sumbangan kajian mengenai masalah pembelajaran ini.*

# TABLE OF CONTENTS

# LIST OF FIGURES

Page

UMS

UNIVERSITI MALAYSIA SABAH

# LIST OF TABLES

# ABBREVATION

| | |
|---|---|
| **1NN** | 1-Nearest Neighbor |
| **ABC** | Adaptive Base Class |
| **ABKD** | Agent Based Knowledge Discovery |
| **ACM** | Association for Computing Machinery |
| **AdaBoost** | Adaptive Boosting |
| **AdaBoostM1** | Adaptive Boosting Multiclass1 |
| **Adacost** | Adaptive Boosting with Cost (Misclassification cost-sensitive boosting) |
| **ANN** | Artificial Neural Networks |
| **AUC** | Area Under Curve |
| **AVA** | All-Versus-All |
| **BABoost** | Balanced AdaBoost |
| **BHS** | Binary Hierarchical Classifier |
| **C4.5** | C4.5 is an algorithm used to generate a decision tree developed by Ross Quinlan |
| **C5.0** | Commercial version of C4.5 (algorithm to generate a decision tree) |
| **CATCH** | Canadian Assessment of Tomography for Childhood Head Injury |
| **CNNDM** | Class Nearest Neighbor Distance Matrix |
| **CT** | Computed Tomography |
| **CWW** | Class Confidence Weight |
| **DAG** | Directed Acyclic Graph |
| **DataBoost** | Data Boosting |
| **DataBoost-IM** | Data Sets with Boosting and Data Generation – Imbalance |

| | |
|---|---|
| **DB2** | Divide-by-2 |
| **DBEG** | Distribution-Based Example Generation |
| **DBKNN** | Density-based-EKNN |
| **DDAG** | Decision Directed Acyclic Graph |
| **DECIML** | Direct Ensemble Classifier for Imbalanced Multiclass Learning |
| **DECIMLFS** | Direct Ensemble Classifier for Imbalanced Multiclass Learning with Feature Selection |
| **DECIMLFS.FIG** | Direct Ensemble Classifier for Imbalanced Multiclass Learning with Filter-based feature selection and Information Gain threshold |
| **DECIMLFS.WIG** | Direct Ensemble Classifier for Imbalanced Multiclass Learning with Wrapper-based feature selection and Information Gain threshold |
| **DECIMLFS.WR** | Direct Ensemble Classifier for Imbalanced Multiclass Learning with Wrapper-based feature selection and Random Information Gain threshold |
| **DS** | Decision Stump |
| **DT** | Decision Tree |
| **ECOC** | Error-Correcting Output Code |
| **eKISS** | ensemble Knowledge for Imbalance Sample Sets |
| **EKNN** | Evidence-theory-based-KNN |
| **F-measure** | F1 Score/Balance F-Score (to measure test accuracy) |
| **FN** | False Negative |
| **FP** | False Positive |
| **FS** | Feature Selection |
| **FSMC** | Feature Selection for Minority Class |
| **GA** | Genetic Algorithm |
| **GC** | Generalized Coding |

| | |
|---|---|
| **G-mean** | Geometric mean |
| **HSVM** | Hierarchical Support Vector Machines |
| **IB1** | Instance Based (learning algorithm in Weka) |
| **ID3** | Iterative Dichotomiser 3 - is an algorithm used to generate a decision tree invented by Ross Quinlan |
| **IEEE** | Institute of Electrical and Electronics Engineers |
| **IG** | Information Gain |
| **IR** | Imbalance Ratio |
| **J48** | Open source Java implementation of the C4.5 algorithm in Weka |
| **KDD** | Knowledge Discovery in Databases |
| **KEEL** | Knowledge Extraction based on Evolutionary Learning (data repository) |
| **KNN** | k-Nearest Neighbor |
| **LI** | Lack of Information |
| **LogitBoost** | Logistic Boosting |
| **LogitBoost-J** | Logistic Boosting (extended for unbalanced data situation) |
| **M1** | Model 1 |
| **M1v** | Model 1 vote |
| **M2** | Model 2 |
| **M2v** | Model 2 vote |
| **MAP** | Maximum A-Posteriori |
| **MDLP** | Minimum Description Length Principle |
| **MGM** | Maximum Geometry Mean |
| **MLNN** | Multiclass Leveraged k-Nearest Neighbor |
| **MLP** | Multi Layer Perceptron |
| **MLP** | Multilayer Perceptron |
| **MMC** | Moving Median Center hypersphere |

| | |
|---|---|
| **MPEG-7** | Multimedia Content Description Interface – 7 |
| **MS** | Maximum Sum |
| **NB** | Naïve Bayes |
| **NIPS** | Neural Information Processing Systems Conference (data repository for feature selection challenge) |
| **OAA** | One-Against-All |
| **OAO** | One-Against-One |
| **OR** | Or truth |
| **OVA** | One-Versus-All |
| **PAC** | Probably Approximately Correct |
| **PAQ** | P-Against-Q |
| **PART** | Projective Adaptive Resonance Theory |
| **ROC** | Receiver Operating Characteristic |
| **RSM** | Random Subspace Method |
| **RUSBoost** | Random Under-Sampling with Boosting |
| **S2N** | Signal to Noise correlation coefficient |
| **SLIPPER** | Simple Learner with Iterative Pruning to Produce Error Reduction |
| **SMO** | Sequential Minimal Optimization (algorithm) |
| **SMOTE** | Synthetic Minority Over-sampling Technique |
| **SMOTEBoost** | Synthetic Minority Over-sampling Technique with Boosting |
| **SVM** | Support Vector Machines |
| **TN** | True Negative |
| **TP** | True Positive |
| **TPR** | True Positive Rate |
| **UCI** | University of California, Irvine |

| | |
|---|---|
| **UCLA** | University of California, Los Angeles |
| **UCR** | University of California, Riverside (data repository) |
| **Weka** | Waikato Environment for Knowledge Analysis |
| **WINNOW** | Machine learning algorithm for learning a linear classifier from labeled examples similar to the perceptron algorithm |
| **fMRI** | functional Magnetic Resonance Imaging |

# SYMBOL

| | |
|---|---|
| $\subseteq$ | Is contained in |
| $\mathbb{R}^n$ | Set of $n$ real number |
| $\mathfrak{R}^n$ | Set of $n$ real number |
| $\mathbb{H}$ | The (set of) quaternions/hypotheses |
| $\in$ | Is an element of |
| $\parallel$ | Parallel; Is parallel to |
| $\mid$ | Conditional probability; given |
| $\Pi$ | Product; product of all values in range of series |
| $\sqrt{x}$ | Square root |
| $\Sigma$ | Summation; sum of all values in range of series |
| $\hat{f}(x)$ | Circumflex/estimator; of the function of $x$ |
| $\delta(x)$ | Delta function |
| $\delta(x) = \begin{cases} x \\ y \end{cases}$ | Dirac delta function; hyperfunction |
| $=$ | Is equal to |
| $==$ | Equivalence |
| $\geq$ | Greater or equal |
| $>$ | Greater than |
| $\beta$ | Beta |
| $\mu$ | Mu |
| $\leftarrow$ | Arrow function; from..to; set theory |
| $\theta$ | Theta |
| $\Delta$ | Symmetric difference |
| $\varepsilon$ | Epsilon; represent small number near zero |

# CHAPTER 1

# INTRODUCTION

## 1.1    Background

Over the decades, knowledge discovery through data mining and machine learning have been extensively studied and applied in various fields. It continues to solve many real world problems and applications, such as pattern recognition, computer vision, image processing, bioinformatics, and a lot more complicated domains. Researches in this massive artificial intelligence domain bring notable advantages where data mining and machine learning algorithms provide assistance in helping people for the knowledge discovery in databases (KDD).

The classic definition of KDD is "the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data" (Fayyad *et al.*, 1996). Knowledge discovery is one of artificial intelligence contributions to research community which includes the processes of gathering data and information, pre-processing, analyzing data, extracting hidden knowledge (with data mining), constructing knowledge, knowledge evaluations and knowledge reuse. Knowledge discovery is not only using artificial intelligence techniques but it is also supported by other theories in database, retrieval theories, computing, visualization, statistics, etc. Generally, the data mining processes are grouped into three major key steps: preparation of input data, mining of data, and post-processing of output patterns (Du, 2010).

Knowledge and decisions applied today were learned from past experiences and they are iteratively refined to be applied in future problems encountered. Data mining is one of the steps in the knowledge discovery process that can be used to automate the discovery of patterns or data modeling for selected empirical data, visualize and finally use the learned knowledge in response to future unseen data. Data mining is defined as the extraction of implicit, previously unknown, and potentially useful knowledge from data (Witten and Frank, 2000). In other words,

data mining can be interpreted as the task of employing an algorithm that processes raw data automatically or semi-automatic and extracts any meaningful patterns that will be used for prediction tasks on new unseen data. Machine learning specifically provides various methods and learning algorithms that can be used to find and describe any structural patterns in data. Thus, the study of machine learning algorithms has emerged as the technical basis for any data mining works.

Currently, many of these common algorithms and their advanced variations provide high classification performance in various empirical data. Advanced techniques, such as ensemble learning methods, are employed which apply different learning algorithms for different applications and these methods provide even higher classification accuracy. For example, ensemble learning approaches have been applied in a shape classification task is able to classify MPEG-7 shape and Swedish leaf shape dataset as high as 95 percent and 98 percent (Temlyakov et al., 2010). As a result, it is substantially harder for new researchers propose any better classification methods. The result of high performance accuracies produced by these advanced techniques indicates that the advancement of data mining methods and machine learning algorithms are almost stagnant and it can solve almost any classification tasks. However, real world problem is far from the reality of no unsolvable classification problem. There is still exist diverse problem requiring different efforts to find efficient solution for comparison, such as new dataset (even in the similar domain which was solved before), massive data, incomplete data (caused by noise or missing values), etc. For example, one of the most challenging data mining problems that are still receiving attention among researchers is the multiclass and imbalance classification problem (Alejo et al., 2008; Ghanem et al., 2010; Lerteerawong and Athimethphat, 2011; Tahir et al., 2010; Valizadegan et al., 2008; Zhou and Liu, 2010).

In data mining, multiclass classification problem refers to assigning one of the several class labels to an input object. Unlike the binary classification, learning a multiclass problem is a more complex task due to the fact that each example can only be assigned to exactly one class label (Valizadegan et al., 2008). In fact, numerous attempts of using binary classification methods have failed to perform

2