

# Case-Base Retrieval of Childhood Leukaemia Patients Using Gene Expression Data

A Thesis Submitted for the Degree of  
Doctor of Philosophy

By

*Ali Anaissi*

in

School of Software  
UNIVERSITY OF TECHNOLOGY, SYDNEY  
AUSTRALIA  
JANUARY 2013

© Copyright by Ali Anaissi, 2013

UNIVERSITY OF TECHNOLOGY, SYDNEY  
SCHOOL OF SOFTWARE

The undersigned hereby certify that they have read this thesis entitled "**Case-Base Retrieval of Childhood Leukaemia Patients Using Gene Expression Data**" by **Ali Anaissi** and that in their opinions it is fully adequate, in scope and in quality, as a thesis for the degree of **Doctor of Philosophy**.

Dated: January 2013

Research Supervisor: \_\_\_\_\_  
Dr Madhu Goyal

# CERTIFICATE

Date: **January 2013**

Author: **Ali Anaissi**

Title: **Case-Base Retrieval of Childhood Leukaemia  
Patients Using Gene Expression Data**

Degree: **Ph.D.**

I certify that this thesis has not already been submitted for any degree and is not being submitted as part of candidature for any other degree.

I also certify that the thesis has been written by me and that any help that I have received in preparing this thesis, and all sources used, have been acknowledged in this thesis.

---

Signature of Author

## Acknowledgements

My first acknowledgement is to God for being my source of spiritual strength. I would like to thank my supervisor Dr Madhu Goyal and Dr Paul Kennedy for the opportunity to work with them on this thesis, for all of the guidance they have given me along the way. I would like to thanks Prof Jie Lu for accepting me to be a member of QCIS lab. The University of Technology Sydney has provided me not only with a wonderful education but also the scholarship that made this work possible and for that I am grateful. The Children's Hospital at Westmead for providing me the gene expression datasets. I thank my friends and colleagues at the QCIS lab for the unforgettable memories. Thank you all.

*To my father and mother who will always be an inspiration throughout my life. To my wife for being my partner in life. And, to my children Joelle and Jawad, who will always be my blessings from God. You two represent my greatest commitments in life. I pray for wisdom so I can accomplish my job as your father. All together, we have survived throughout this, in our situation, challenging journey. I love you all.*

# Table of Contents

<b>Table of Contents</b>	<b>vii</b>
<b>List of Tables</b>	<b>viii</b>
<b>List of Figures</b>	<b>x</b>
<b>Abstract</b>	<b>1</b>
<b>1 Introduction</b>	<b>4</b>
1.1 Motivation . . . . .	6
1.2 Objective and Aims . . . . .	7
1.3 Contributions . . . . .	8
1.4 Organization of this Thesis . . . . .	10
<b>2 Literature Review</b>	<b>11</b>
2.1 Case-Based Reasoning . . . . .	11
2.1.1 Case-Based Retrieval For Gene Expression Data . . . . .	13
2.1.2 $k$ -Nearest Neighbour Classifier . . . . .	15
2.2 Feature Selection . . . . .	17
2.2.1 Feature Selection Techniques . . . . .	18
2.2.2 Adaptation of Ensemble Methods in Feature Selection . . . . .	21
2.2.3 Random Forest For Feature Selection . . . . .	22
2.3 Dimensionality Reduction . . . . .	24
2.3.1 Linear Dimensionality Reduction . . . . .	24
2.3.2 Non-Linear Dimensionality Reduction . . . . .	26
2.3.3 Trustworthiness . . . . .	27
2.4 Genetic Algorithms . . . . .	29
2.5 Imbalanced Data . . . . .	31
2.5.1 Approaches to Handle Imbalanced Classes . . . . .	31

2.5.2	Sampling Technique . . . . .	32
2.5.3	Cost Sensitive Learning . . . . .	33
2.6	Evaluation Measures of Imbalanced Classes . . . . .	33
2.7	Summary . . . . .	35
<b>3</b>	<b>Case-Base Retrieval Framework For Gene Expression Datasets</b>	<b>37</b>
3.1	Case Base Retrieval in Gene Expression Datasets . . . . .	37
3.2	Case-Base Retrieval Framework . . . . .	39
3.2.1	Module 1: Pre-processing the Training Dataset . . . . .	41
3.2.2	Module 2: Pre-processing the test dataset . . . . .	45
3.3	Summary . . . . .	47
<b>4</b>	<b>Balanced Iterative Random Forest for Gene Selection from Microarray Data</b>	<b>48</b>
4.1	Balanced Iterative Random Forest for Feature Selection . . . . .	49
4.1.1	Handling Imbalanced Classes Effect on Feature Selection . . . . .	50
4.1.2	Algorithm of Balanced Iterative Random Forest . . . . .	54
4.1.3	Validation of Over-Fitting . . . . .	55
4.1.4	Validation of the Selected Genes . . . . .	56
4.2	Datasets . . . . .	57
4.3	Experiments . . . . .	60
4.3.1	Experiments on Childhood Leukaemia Dataset . . . . .	60
4.3.2	Experiments on the Three Public Microarray Datasets . . . . .	64
4.4	Comparison With Other Algorithms . . . . .	66
4.5	Summary . . . . .	67
<b>5</b>	<b>Dimensionality Reduction and Visualization</b>	<b>69</b>
5.1	Dimensionality Reduction of Gene Expression Datasets . . . . .	70
5.1.1	Linear Dimensionality Reduction Approach . . . . .	70
5.1.2	Non-Linear Dimensionality Reduction and Visualization of Gene Expression Datasets . . . . .	72
5.2	Experiments . . . . .	74
5.2.1	Choosing number of Components in KPCA . . . . .	75
5.2.2	Performance Evaluation of LPC . . . . .	77
5.3	Summary . . . . .	86
<b>6</b>	<b>Similarity Measurement</b>	<b>90</b>
6.1	Similarity Measurement Approach . . . . .	90
6.2	Experiments . . . . .	92

6.2.1	Distance Metric Selection . . . . .	92
6.2.2	Determination of Parameter $k$ . . . . .	94
6.2.3	Feature Weighting . . . . .	95
6.2.4	Oversampling . . . . .	101
6.3	Summary . . . . .	104
<b>7</b>	<b>Conclusion</b>	<b>109</b>
7.1	Contribution 1: Case-Base Retrieval Framework . . . . .	110
7.2	Contribution 2: Balanced Iterative Random Forest for Feature Selection . . . . .	110
7.3	Contribution 3: Local Principal Component Analysis for Dimensionality Reduction and Visualisation . . . . .	111
7.4	Contribution 4: Weight Learning Genetic Algorithm for Feature Weighting . . . . .	112
7.5	Future Work . . . . .	113
	<b>Bibliography</b>	<b>116</b>

# List of Tables

2.1	A confusion matrix for a two-class classification . . . . .	34
4.1	Microarray gene expression datasets . . . . .	59
4.2	A confusion matrix for the childhood leukaemia training dataset . . .	61
4.3	A confusion matrix for the childhood leukaemia test dataset . . . . .	61
4.4	A confusion matrix for the childhood leukaemia test dataset (first list)	64
4.5	A confusion matrix for the childhood leukaemia test dataset (second list)	64
4.6	A confusion matrix for the childhood leukaemia test dataset (third list)	64
4.7	A confusion matrix for the Golub training dataset . . . . .	65
4.8	A confusion matrix for the Golub test dataset . . . . .	65
4.9	A confusion matrix for the Colon training dataset . . . . .	66
4.10	A confusion matrix for the Colon test dataset . . . . .	66
4.11	A confusion matrix for the Lung cancer training dataset . . . . .	67
4.12	A confusion matrix for the Lung cancer test dataset . . . . .	67
4.13	Accuracy results for Colon and Leukaemia datasets . . . . .	68
6.1	Performance comparison of the distance metrics for the nearest neighbour classifier. . . . .	94
6.2	Performance comparison of the $k$ NN for different values of $k$ . . . . .	95
6.3	Classification performance results of the test dataset . . . . .	95
6.4	Classification performance results of the test data applied on the weighted-5NN classifier . . . . .	100
6.5	Classification probability results of the test dataset . . . . .	106

6.6	Classification performance of the one-side over-sampled . . . . .	107
6.7	Classification performance of 100% over-sampled training dataset . .	107
6.8	Classification performance results of the test dataset after 100% over-sampled the training dataset . . . . .	107
6.9	Classification probability results of the test dataset after SMOTE . .	108

# List of Figures

2.1	Case-base reasoning stages . . . . .	12
2.2	Selection methods. (a) Filter approach (b) Wrapper approach (c) Embedded approach [78] . . . . .	19
2.3	Neighbours problem in dimensionality reduction: a) high dimensional space and b) low dimensional space . . . . .	28
2.4	Genetic algorithm process based on [70] . . . . .	30
3.1	Case-base retrieval framework . . . . .	40
3.2	Pre-processing the training dataset in the framework . . . . .	42
3.3	Modification of the data in the dimensionality reduction and visualisation section . . . . .	44
3.4	Modification of the data in the feature weighting and sampling techniques section . . . . .	45
3.5	Pre-processing the test sample in the framework and retrieving the similar cases . . . . .	46
4.1	Variations of Out-of-bag error during selecting the features . . . . .	62
5.1	Structure of the input and output dataset of LPC . . . . .	72
5.2	Accuracy according to the different dimensionality reduction process using 8-fold cross-validation . . . . .	77
5.3	Trustworthiness of LPC on Swiss-roll data set . . . . .	79
5.4	Trustworthiness of LPC on Childhood Leukaemia gene expression dataset	80

5.5	Original Iris data set . . . . .	81
5.6	Iris data set processed by LPC . . . . .	82
5.7	Swiss roll data set . . . . .	83
5.8	Swiss roll data reduced to 2D by PCA . . . . .	84
5.9	Swiss roll data reduced to 2D by LPC . . . . .	85
5.10	Leukaemia data reduced to 2D by PCA . . . . .	86
5.11	Leukaemia data reduced to 2D by LPC . . . . .	87
5.12	Trustworthiness of LPC Vs. LLE using Swiss-roll data set . . . . .	88
5.13	Trustworthiness of LPC Vs. LLE using Microarray data set . . . . .	89
6.1	Weight Learning GA and 5NN . . . . .	98
6.2	Fitness value for each generation . . . . .	105

# Abstract

Acute Lymphoblastic Leukaemia (ALL) is the most common childhood malignancy. Nowadays, ALL is diagnosed by a full blood count and a bone marrow biopsy. With microarray technology, it is becoming more feasible to look at the problem from a genetic point of view and to perform assessment for each patient. This thesis proposes a case-base retrieval framework for ALL using a nearest neighbour classifier that can retrieve previously treated patients based on their gene expression data. However, the wealth of gene expression values being generated by high throughput microarray technologies leads to complex high dimensional datasets, and there is a critical need to apply data-mining and computational intelligence techniques to analyse these datasets efficiently. Gene expression datasets are typically noisy and have very high dimensionality. Moreover, gene expression microarray datasets often consist of a limited number of observations relative to the large number of gene expression values (thousands of genes). These characteristics adversely affect the analysis of microarray datasets and pose a challenge for building an efficient gene-based similarity model. Four problems are associated with calculating the similarity between cancer patients on the basis of their gene expression data: feature selection, dimensionality reduction, feature weighting and imbalanced classes.

The main contributions of this thesis are: (i) a case-base retrieval framework, (ii) a

Balanced Iterative Random Forest algorithm for feature selection, (iii) a Local Principal Component algorithm for dimensionality reduction and visualization and (iv) a Weight Learning Genetic algorithm for feature weighting.

This thesis introduces Balanced Iterative Random Forest (BIRF) algorithm for selecting the most relevant features to the disease and discarding the non-relevant genes. Balanced iterative random forest is applied on four cancer microarray datasets: Childhood Leukaemia dataset, Golub Leukaemia dataset, Colon dataset and Lung cancer dataset. Childhood Leukaemia dataset represents the main target of this project and it is collected from The Children's Hospital at Westmead. Patients are classified based on the cancer's risk type (Medium, Standard and High risk); Colon cancer (cancer vs. normal); Golub Leukaemia dataset (acute lymphoblastic leukaemia vs. acute myeloid leukaemia) and Lung cancer (malignant pleural mesothelioma or adenocarcinoma). The results obtained by BIRF are compared to those of Support Vector Machine-Recursive Feature Elimination (SVM-RFE) and Naive Bayes (NB) classifiers. The BIRF approach results are competitive with these state-of-art methods and better in some cases. The Local Principal Component (LPC) algorithm introduced in this thesis for visualization is validated on three datasets: Childhood Leukaemia, Swiss-roll and Iris datasets. Significant results are achieved with LPC algorithm in comparison to other methods including local linear embedding and principal component analysis. This thesis introduces a Weight Learning Genetic algorithm based on genetic algorithms for feature weighting in the nearest neighbour classifier. The results show that a weighted nearest neighbour classifier with weights generated from the Weight Learning Genetic algorithm produces better results than the un-weighted nearest neighbour algorithm.

This thesis also applies synthetic minority over sampling technique (SMOTE) to increase the number of points in the minority classes and reduce the effect of imbalanced classes. The results show that the minority class becomes recognised by the nearest neighbour classifier. SMOTE also reduces the effect of imbalanced classes in predicting the class of new queries especially if the query sample should be classified to the minority class.