

# A MODEL BASED APPROACH FOR CLUSTER TRADITIONAL RICE VARIETIES OF SRI LANKA

Mapulage Don Rangika Lilukshi Silva

08/8022



University of Moratuwa, Sri Lanka.  
Electronic Theses & Dissertations  
Degree of Master of Philosophy

Department of Computer Science and Engineering

University of Moratuwa  
Sri Lanka

January 2015

## Declaration of the candidate & Supervisor

I declare that this is my own work and this thesis does not incorporate without acknowledgment any material previously submitted for a Degree or Diploma in any other University or Institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgment is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature: .....

Date:.....



University of Moratuwa, Sri Lanka.  
Electronic Theses & Dissertations  
[www.lib.mrt.ac.lk](http://www.lib.mrt.ac.lk)

The above candidate has carried out research for the MPhil thesis under my supervision.

Signature of the supervisor:.....

Date:.....

## ACKNOWLEDGMENT

First and foremost, I would like to express my special thanks and appreciation to my supervisor Professor Nalin Wickramarachchi, for his guidance, advice, knowledge and word of encouragement during this study. I am very much grateful to him, for identifying my capabilities and letting me step into the world of Bioinformatics which opened a great research path for me, without whom it would not have been possible. His advice on both research as well as on my career have been priceless and his kindness will always be remembered.

I would like to extend my thanks with gratitude to the former head of the department, Computer Science and Engineering, Ms. Vishaka Nanayakkara and Lt.Col. Dr. Chandana Gamage, former M.Sc. coordinator associate professor Sanath Jayasena and present M.Sc. coordinator Dr. Shehan Perera for the opportunity given to pursue my post graduate studies at the department and for the immense support and guidance given throughout my period. Also, I would like to thank the entire faculty at the department of Computer Science and Engineering for making my stay at the department a memorable experience.

My sincere appreciation also goes to Dr. Jagath Weerasinghe and Dr. Gowri Rajkumar at Institute of Biochemistry, Molecular Biology and Biotechnology, University of Colombo for providing me the biological data needed to evaluate my research methodology. It is a great honor to be a member of their research team and the opportunity given was a great help to acquire the biological knowledge needed to carry out my research work.

I must acknowledge the immeasurable contributions of my friends and colleagues who have shown great love and care during my study.

This acknowledgment won't be complete without my family. My husband Saman has been the greatest strength behind my success and his support made me to achieve many goals in my life. I am thankful for his patient and support given throughout my life. A profound gratitude goes to my mother and to my two sisters whom gave me a great support by looking after my three children in my absence.

## ABSTRACT

As a result of the enormous volume of data produced by highly developed modern techniques, focus on clustering biological data has shown a great interest among biologist to detect the underlying patterns in data since the biological experiment itself has failed to identify the hidden information and divergence patterns exist in data correctly.

This study aims to (1) assist clustering biologically similar sequences to detect divergence patterns exist in rice genomic data, by developing a program using the model based clustering algorithm based on Chinese restaurant process which was originally proposed to cluster gene expression data (2) focus on finding the performance of calculating the pairwise distance matrix of rice genome sequences based on the 12-dimensional natural vector of the DNA sequence, as the similarity measure in cluster analysis.

The developed program based on the proposed model based clustering method was executed on ALFP profile data set consisting features of 53 Sri Lankan traditional and wild rice varieties in order to identify the genetic divergence among them. Both a statistical and a biological cluster evaluation were carried out to validate the results obtained. Statistical evaluation was done based on the Bayes ratio to measure the tightness of the clusters formed. Biological evaluation was conducted with the help of the domain experts and research work done by the institute of rice of Sri Lanka. The results showed that the proposed algorithm is capable of identifying highly similar varieties of rice showing their divergence patterns.

Finding the performance of how well the natural vector method captures the information encoded in rice genome sequences, 10 rice disease resistance genes which belong to three different protein families from Rice genome annotation project database were used. The results showed that the pairwise distance matrix calculated based on 12-dimensional natural vector method gives efficient results compared to traditional proximity matrices. It also revealed that the fixed length size sequences (sub sequences) which are not greater than the minimum total length of the selected sequences are also highly capable of capturing the encoded information in total length, regardless of the sub sequence length.

**Key words: Model-Based clustering, Genetic Diversity**

## Table of contents

<b>Declaration of the candidate &amp; Supervisor</b>	<b>i</b>
<b>Acknowledgement</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>Table of Contents</b>	<b>iv</b>
<b>List of figures</b>	<b>vi</b>
<b>List of tables</b>	<b>vii</b>
<b>List of abbreviations</b>	<b>viii</b>
<b>1 INTRODUCTION</b>	<b>1</b>
1.1 Problem Statement	2
1.2 Aim of This Work	3
1.3 Approach	3
1.4 Scope and Limitation	4
1.5 Thesis Outline	4
<b>2 BIOLOGICAL BACKGROUND</b>	<b>6</b>
2.1 The Molecular Building Blocks of Life	6
2.1.1 Deoxyribonucleic Acid (DNA)	8
2.1.2 Ribonucleic Acid (RNA)	9
2.1.3 Proteins	10
2.1.4 Central dogma of molecular biology	12
2.2 Regulation of Genes	14
<b>3 STATISTICAL BACKGROUND</b>	<b>15</b>
3.1 Knowledge Discovery and Data Mining	15
3.2 Machine Learning	16
3.3 Cluster Analysis	17
3.3.1 Cluster definition	18
3.3.2 Clustering algorithms	21

---

3.3.3	Other Approaches . . . . .	25
<b>4</b>	<b>LITERATURE REVIEW</b>	<b>30</b>
4.1	Previous Work . . . . .	30
<b>5</b>	<b>METHODOLOGY</b>	<b>33</b>
5.1	Statistical Model . . . . .	34
5.1.1	Predictive update . . . . .	36
5.2	Clustering Algorithm . . . . .	38
5.2.1	Proximity measurements for DNA sequences . . . . .	39
5.2.2	Assessing the suitability . . . . .	40
<b>6</b>	<b>RESULTS AND DISCUSSION</b>	<b>42</b>
6.1	Data Sets . . . . .	42
6.1.1	Sri Lankan Rice varieties data set . . . . .	42
6.1.2	Oryza sativa data set . . . . .	43
6.2	Experimental Design . . . . .	44
6.2.1	Sri Lankan rice varieties data analysis . . . . .	44
6.2.2	DNA sequence analysis - Oryza sativa data set . . . . .	47
6.3	Discussion . . . . .	54
<b>7</b>	<b>CONCLUSION AND FUTURE WORK</b>	<b>57</b>
7.1	Conclusions . . . . .	57
7.2	Future Work . . . . .	58
	<b>References</b>	<b>60</b>
	<b>Appendix A: Input Data sets</b>	<b>69</b>



University of Moratuwa, Sri Lanka.  
 Electronic Theses & Dissertations  
[www.lib.mrt.ac.lk](http://www.lib.mrt.ac.lk)

## LIST OF FIGURES

	<b>Page</b>
2.1 The molecule of life [1] . . . . .	7
2.2 DNA and its building blocks [2] . . . . .	9
2.3 Primary protein structure [3] . . . . .	11
2.4 Central Dogma of Molecular Biology [4] . . . . .	13
3.1 KDD process [5] . . . . .	16
3.2 Example of clustering . . . . .	18
3.3 Example of dendrogram . . . . .	23
3.4 Density based data sets [6] . . . . .	25
3.5 Kohenan Self-Organization Map . . . . .	27
6.1 Dendrogram drawn for whole length sequences with natural vector method . . . . .	49
6.2 Dendrogram drawn for whole length sequences with Clustal Omega program . . . . .	50
6.3 Dendrogram drawn for first 600bp with natural vector method . .	51
6.4 Dendrogram drawn for first 700bp with natural vector method . .	52
6.5 Dendrogram drawn for first 800bp with natural vector method . .	52
6.6 Dendrogram drawn for first 900bp with natural vector method . .	53
6.7 Dendrogram drawn for first 800bp with Clustal Omega program .	53
1 The UPGMA dendrogram showing genetic diversity among Sri Lankan traditional rice varieties . . . . .	71

## LIST OF TABLES

	<b>Page</b>
2.1 The twenty amino acids and their abbreviation . . . . .	12
3.1 similarity measures for continuous features . . . . .	20
3.2 Similarity measures for binary data . . . . .	20
6.1 Protein super families and selected gene sequences with their length in base pairs. . . . .	43
6.2 Clustering results . . . . .	45
6.3 Euclidean pairwise distance matrix obtained by natural vectors for whole length sequences . . . . .	49
1 Jaccard's pairwise distances for the 53 rice varieties . . . . .	70



University of Moratuwa, Sri Lanka.  
Electronic Theses & Dissertations  
[www.lib.mrt.ac.lk](http://www.lib.mrt.ac.lk)