AIN SHAMS UNIVERSITY



Faculty of Computer & Information Sciences *Abbassia, Cairo, Egypt*

Data Mining Techniques in Gene Expressions

A Thesis submitted to the Department of Computer Science, Faculty of Computer and Information Sciences, Ain Shams University,

In partial fulfillment of the requirements for the degree of Master of Computer Science

By

Basma Ali Maher

B.Sc. in Computer Science, 2009 Computer Science Department, Faculty of Computer and Information Sciences, Zagazig University

Under the Supervision of

Prof. Dr. Abd El-Badeeh M. Salem

Professor of Computer Science, Computer Science Department, Faculty of Computer and Information Sciences, Ain shams University

Prof. Dr. El-Sayed Mohamed El-Horbaty

Professor and head of Computer Science Department, Faculty of Computer and Information Sciences, Ain shams University

Dr. Abeer Mahmoud Mahmoud

Lecturer,

Computer Science Department,
Faculty of Computer and Information Sciences,
Ain shams University

Acknowledgement

First and foremost, I would like to thank *Allah* for giving me the opportunity and the strength to accomplish this work.

I would like to express my greatest and deepest gratitude to **Prof. Dr. Abd El-Badeeh Mohammed Salem and Prof. Dr. El-Sayed Mohamed El-Horbaty** for planning and full supervision of this work, continuous advises, guidance and greatest help in interpretation of the results. Thanks a lot for their valuable suggestions in writing the thesis.

A very deep and special thanks to my direct supervisor **Dr. Abeer Mahmoud Mahmoud** for her continuous help, time and effort during the research years of this study.

I would like to express my deep grateful thanks to **My Family** who encouraged me too much during the progress of the work.

Abstract

In recent years, the rapid developments in the genetics field have generated a huge amount of biological data. Microarray gene expression data is an important instance of biological data. It has high dimensionality with a small number of samples accompanied with large number of genes. Therefore, using machine learning techniques for knowledge discovery in such data become a rich area for researchers. The mining phase is usually divided into two steps: the gene selection (feature reduction) and the classification process.

Gene selection is a process of finding the genes most strongly related to a particular class. The benefit of this process is to reduce not only dimensionality but also, the danger of presence of irrelevant genes that affect the classification process. Many machine learning approaches are used feature reduction but the study focuses on t-test and class separability. In the other hand, classification is an important data-mining problem that has a wide range of applications. Classification concerns learning that classifies data into the predetermined categories. It is applied to discriminate diseases or to predict outcomes based on gene expression patterns and perhaps even identify the best treatment for given genetic signature. Many machine learning approaches are used classification. In this study, it focuses on Support vector machine and k-nearest neighbor.

Support Vector Machine (SVM) plays a very important role in the data mining classification problem. The structure of SVM depends on kernel functions, where the most commonly used are liner and polynomial. If there are more than two classes in the data set, binary SVMs are not sufficient to solve the whole problem. To solve multi-class classification problems, the whole problem should be converted into a number of binary classification problems. Usually, there are two approaches. One is the "one against all" scheme and the other is the "one against one" scheme.

On the other hand, K-Nearest Neighbor shows an outstanding performance in many cases of classifying microarray gene expression. For using KNN technique three key elements are essential, (1) a set of data for training, (2) a group of labels for the training data (identifying the class of each data entry) and (3) the value of K for deciding the number of nearest neighbors.

This study proposes a new hybrid reduction approach for the promotion of the cancer classification accuracy that uses two gene selection techniques to confirm the most informative genes and to discard irrelevant genes that affect the classification accuracy. Actually, it applied two machine learning (ML) gene ranking techniques (T-test and Class Separability (CS)) and two ML classifiers; K-nearest neighbor (KNN) and support vector machine (SVM); for exploring and analyzing the process of mining microarray gene expression profiles. In addition, based on these analyses we proposed a hybrid ML reduction approach to enhance the classification accuracy.

It has tested and validated the ML approaches on four public microarray databases; Lymphoma, Leukemia, Small Round Blue Cell Tumors (SRBCT) and Lung Cancer datasets. The experimental results show that the hybrid system achieves enhancement in the classification accuracy better than the SVM and KNN techniques alone. Also, selecting genes from the whole data is better than selecting it from the training data. But excluding the testing samples from the classifier building process, make it more accurately to compare the performance and it make a validation for the system.

List of Publications

- [1]. Basma A.Maher, Abeer M. Mohamed, El-Sayed M.El-horbaty and Abd El-Badeeh M.Salem, "Classification of Two Types of Cancer Based on Microarray Data", Egyptian Computer Science (ECS), Vol. 38, No.2, PP.56-66, 2014.
- [2]. Abeer M. Mohamed, Basma A.Maher, "A Hybrid Reduction Approach for Enhancing Cancer Classification of Microarray Data", Int. Journal of Advanced Research in Artificial Intelligence, Vol. 3, No. 10, PP. 1-10, 2014.
- [3]. Abeer M. Mohamed, Basma A.Maher, El-Sayed M.El-Horbaty and Abd El-Badeeh M. Salem, "Applying a Statistical Technique for the Discovery of Differentially Expressed Genes in Microarray Data", Proceedings of the 4th European Conference of Systems (ECS '13), France, PP. 220-227, 2013.
- [4]. Abeer M. Mohamed, Basma A.Maher, El-Sayed M.El-Horbaty and Abd El-Badeeh M. Salem. "Analysis of Machine Learning Techniques for Gene Selection and Classification of Microarray Data", Proceedings of the 6th Int. Conf. on Information Technology, Cloud Computing, Jordon, 2013.

Table of Contents

Ack	nowled	gement	I
Abs	tract		II
List	of Pub	lications	IV
List	of Figu	ires	VII
	of Tab		VIII
List	of Abb	reviations	XI
1-	Introd	luction	1
	1.1 1.2 1.3 1.4 1.5 1.6	Problem Definition Research Motivation Thesis Objectives Contribution Methodology Thesis Organization	1 3 3 4 5 5
2-	Litera	ture Review and Related Works	7
	2.1 2.2 2.3 2.4	Data Mining Research on Bioinformatics 2.2.1 Microarray Gene Expression. 2.2.2 Gene Selection Techniques. 2.2.3 Gene Classification Techniques. Related Works. Summary.	8 9 10 10 11 11
3-	Gene E	Expression Microarray Databases	24
	3.1 3.2 3.3 3.4 3.5	Introduction Lymphoma Dataset. Leukemia Dataset. SRBCT Dataset. Lung Cancer Dataset. Summary	24 26 28 29 30

4-	Ma	chine Learning Techniques for the Discovery of DEGs
	4.1	Introduction
	4.2	Gene Selection Techniques
	4.3	Considering All Data Reduction
		4.3.1 T-test
		4.3.2 Class Separability
	4.4	Considering Train Data only
		4.4.1 T-test
		4.4.2 Class Separability
	4.5	Summary
5-	Ma	chine Learning Classification Approaches of Microarray Data.
	5.1	Introduction
	5.2	Gene Classification Techniques.
		5.2.1 Linear Discriminate Analysis (LDA)
		5.2.2 K-Nearest Neighbor (KNN)
		5.2.3 Support Vector Machine (SVM)
		5.2.4 Fuzzy Neural Network (FNN)
	5.3	Considering All Data Reduction
		5.3.1 K-Nearest Neighbor.
	- 1	5.3.2 Support Vector Machine
	5.4	Considering Train Data only
		5.4.1 K-Nearest Neighbor
	<i></i>	5.4.2 Support Vector Machine
	5.5 5.6	Comparisons and DiscussionsSummary
6-	A	Hybrid Reduction Technique for Enhancing Cancer
	Cla	ssification of Microarray Data
	6.1	Mining Workflow
	6.2	Gene Ranking Common List.
	6.3	Experiments and Results
	6.4	Summary
7-	Cor	nclusion and Future Work
	Ref	erences
		ssarv

List of Figures

Figure 3.1	Expression Data Matrix	24
Figure 5.1	The Structure of The FNN	51
Figure 5.2	The Proposed Scenario for Achieving Accurate Classification Accuracy in Microarray	52
Figure 5.3	The Testing Classification Accuracy for The Three Cases of Lymphoma Using KNN	53
Figure 5.4	The Leukemia Testing Classification Accuracy Using T-test with KNN	54
Figure 5.5	The Leukemia Testing Classification Accuracy Using CS with KNN	55
Figure 5.6	The SRBCT Testing Classification Accuracy Using T-test with KNN	56
Figure 5.7	The Leukemia Testing Classification Accuracy Using T-test with SVM	59
Figure 5.8	The Leukemia Testing Classification Accuracy Using CS with SVM	60
Figure 5.9	Proposed Scenario for Gene Expression Data Mining Workflow	61
Figure 5.10	The Testing Classification Accuracy of KNN for The Three Cases of Lymphoma	63
Figure 5.11	Classification Accuracy of KNN and SVM on Lymphoma Three Cases	68
Figure 5.12	The Classification Accuracy of T-test and CS for Both KNN and SVM on Leukemia.	68
Figure 5.13	Six Binary SVM Classifiers with KNN Classifier on SRBCT Database	70
Figure 5.14	Classification Accuracy of KNN and SVM Classifiers on SRBCT	70
Figure 5.15	Classification Accuracy of KNN and SVM for Both T-test and CS on Lung Cancer	72
Figure 6.1	The Proposed Hybrid Reduction Approach of Microarray Data	80

List of Tables

Table	2.1	Comparative Study Analysis	20
Table	2.2	A Comparison of Three Gene Expression Datasets with	22
11		Different Techniques	
Table		A Sample Data from Lymphoma Dataset	27
Table	3.2	Lymphoma Dataset with Empty Spots	27
Table	3.3	Pre-processed Lymphoma Dataset.	27
Table	3.4	The Percentage of Training and Testing Sample for The Three Cases for Lymphoma Dataset	28
Table	3.5	A Sample Data from Leukemia Dataset	29
Table	3.6	A Sample Data from SRBCT Dataset	29
Table	3.7	A Sample Data from Lung Cancer Dataset	30
Table	4.1	Informative Genes of Lymphoma Based on Their T-Test	38
Table	4.2	A Comparison of Our T-test Result with Another Two Results for Lymphoma.	39
Table	4.3	Informative Genes Based on Their T-test for Leukemia Dataset	39
Table	4.4	Informative Genes Based on Their T-test for SRBCT Dataset	40
Table	4.5	A Comparison of Our T-test Result with Another One for The SRBCT Dataset	41
Table	4.6	Informative Genes Based on Their CS for Leukemia Dataset	41
Table	4.7	Top 15 Prioritized Genes and Their Corresponding T-test for The Lymphoma	42
Table	4.8	Leukemia Prioritized Genes Based on Their T-test of Training Dataset	43
Table	4.9	Informative Genes Based on Their T-test for Training SRBCT Dataset.	43
Table	4.10	Lung Cancer Prioritized Genes Based on Their T-test of Training Dataset	44
Table	4.11	Leukemia Prioritized Genes Based on Their CS of Training Dataset	45
Table	4.12	Lung Cancer Prioritized Genes Based on Their CS of Training Dataset	45
Table	5.1	The Testing Classification Accuracy for The Three Cases of Lymphoma Using KNN	53

Table	5.2	The Leukemia Testing Classification Accuracy Using T-test with KNN	54
Table	5.3	The Leukemia Testing Classification Accuracy Using CS with KNN.	55
Table	5.4	The SRBCT Testing Classification Accuracy Using T-test with KNN	56
Table	5.5	The Lymphoma Testing Classification Accuracy for Case 1 Using SVM.	57
Table	5.6	The Lymphoma Testing Classification Accuracy for Case 2 Using SVM	57
Table	5.7	The Lymphoma Testing Classification Accuracy for Case 3 Using SVM	58
Table	5.8	The Leukemia Testing Classification Accuracy Using T-test with SVM	58
Table	5.9	The Leukemia Testing Classification Accuracy Using CS with SVM.	59
Table	5.10	The SRBCT Testing Classification Accuracy Using T-test with SVM.	60
Table	5.11	The Testing Classification Accuracy for The Three Cases of Training Data of Lymphoma Using KNN	62
Table	5.12	The Leukemia Testing Classification Accuracy for Training Using T-test with KNN	63
Table	5.13	The Leukemia Testing Classification Accuracy for Training Using CS with KNN	64
Table	5.14	The SRBCT Testing Classification Accuracy for Training Using KNN	64
Table	5.15	The Lung Cancer Testing Classification Accuracy for Training Using T-test with KNN	65
Table	5.16	The Lung Cancer Testing Classification Accuracy for Training Using CS with KNN	65
Table	5.17	The Lymphoma Testing Classification Accuracy for Training Data for Case 1 Using SVM	66
Table	5.18	The Lymphoma Testing Classification Accuracy for Training Data for Case 2 Using SVM	67
Table	5.19	The Lymphoma Testing Classification Accuracy for Training Data for Case 3 Using SVM	67
Table	5.20	The Leukemia Testing Classification Accuracy for Training Using T-test with SVM	69
Table	5.21	The Leukemia Testing Classification Accuracy for Training Using CS with SVM	69

Table 5.22	The SRBCT Testing Classification Accuracy for Training Using SVM	71
Table 5.23	The Lung Cancer Testing Classification Accuracy for Training Using T-test with SVM	72
Table 5.24	The Lung Cancer Testing Classification Accuracy for Training Using CS with SVM	73
Table 5.25	Comparison of Lymphoma Classifiers	73
Table 5.26	Comparison of Leukemia Classifiers	74
Table 5.27	Comparison of SRBCT Classifiers	75
Table 5.28	Test Classification Accuracy of SVM and KNN Using T-test and CS on Training Data	75
Table 5.29	Comparison of Lymphoma Classifiers	76
Table 5.30	Comparison of Leukemia Classifiers	76
Table 5.31	Results for The SRBCT Dataset Obtained by Different Approaches	77
Table 6.1	The CommonList 24 Genes for Leukemia Cancer and Lung Cancer	80
Table 6.2	Comparison of Leukemia Classifiers	81

List of Abbreviations

ADCA Adenocarcinoma

AI Artificial Intelligence

ALL Acute Lymphoblastic Leukemia

AML Acute Myeloid Leukemia ANN Artificial Neural Network

ANOVA Analysis-of-Variance
BL Burkitt Lymphomas
CC Correlation Coefficient

CLL Chronic Lymphocytic Lymphoma

CS Class Separability
CV Cross Validation

DEGs Differentially Expressed Genes

DLBCL Diffuse Large B-cell Lymphoma

DNA Deoxyribonucleic Acid

E Entropy

ED Euclidean Distance

ELM Extreme Learning Machine
EWS Ewing Families of Tumors
FC Fisher Discriminate Criterion

FL Follicular Lymphoma

FNN Fuzzy Neural Network

FS F(x) Score

KDD Knowledge Discovery Database

KNN K-Nearest Neighbor

LDA Linear Discriminate Analysis

MD Mean Difference
ML Machine Learning

MPM Malignant Pleural Mesothelioma

NB Neuroblastoma

NPPC Nonparallel Plane Proximal Classifier

PCA Principal Component Analysis

RBF Radial Basis Function

RFE Recursive Feature Elimination

RMS Rhabdomyosarcoma

RNA Ribonucleic Acid SANN Subsequent ANN

SNR Signal to Noise Ratio

SRBCT Small Round Blue Cell Tumors

SVM Support Vector Machine

SVM-OAA Support Vector Machine One Against All

TS T-test Statistics

Chapter 1

Introduction

1.1 Problem Definition

Creatures consist of organisms and every organism carries the same genetic information. This genetic information is represented in the form of genes, where only a subset of these genes is active or expressed. Bioinformatics, or computational biology, is the interdisciplinary science of interpreting biological data using information technology and computer science [1]. The broad use of machine learning techniques and their applicability in the different areas of bioinformatics reported a success resolving biological problems because it facilitates the process of analysing such data sets and extracting the important hidden knowledge. The cancer classification based on the microarray data is one example of this type of analysis.

Simply, Microarray gene expression data refers to such repositories of gene information that made the technology of modern biological research. Its goal is to understand the regulatory mechanism that governs protein synthesis and activity of genes. All the cells in an organism carries equal number of genes yet their protein synthesis can be different due to regulation. Protein synthesis is regulated by control mechanisms at different stages [2]:

- 1) Transcription
- 2) RNA splicing
- 3) Translation
- 4) post transitional modifications

Furthermore, analyzing the gene with respect to whether and to what degree they are expressed can help characterize and understand their functions. It can further be analyzed how the activation level of genes changes under different conditions such as for specific diseases (e.g. cancers are generally caused by abnormalities in the genetic material of the transformed cells or change in their activation or function) [3]. Actually, microarray represents a powerful tool in biomedical discoveries and harnessing the potential of this technology depends on the development of appropriate mining approaches [4][5][6].

Microarray techniques provide a plat form where one can measure the expression levels of thousands of genes in hundreds of different conditions. Actually, there is a high redundancy in microarray data and numerous genes contain inappropriate information for precise classification of diseases or phenotypes [7]. Therefore, the amount of data generated by this technology presents a challenge for the biologists to carry out analysis [8].

The mining phase in the knowledge discovery process can be defined as the process of discovering interesting and unknown patterns from large amounts of data stored in information repositories [9][10]. The mining task could be one of regression, summarization, clustering and classification [9]. In microarray data, classification is momentously necessary for cancer diagnosis and treatment. Specially, in classification analysis of microarray data; where the data has high dimensionally; gene selection is one of the critical aspects, where the objective is reaching an efficient gene selection approach that can drastically ease computational burden of the subsequent classification task, and can yield a much smaller and more compact gene set without the loss of classification accuracy.