# A HETEROGENEOUS DATA ENSEMBLE APPROACH FOR PROTEIN FUNCTION PREDICTION UNDER MITOCHONDRION ORGANIZATION

Dinithi Navodhya Sumanaweera

158013D

Thesis/Dissertation submitted in partial fulfillment of the requirements for the degree of Master of Science (Research) in Computer Science

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

October 2016

# DECLARATION

I declare that this is my own work and this thesis/dissertation does not incorporate without acknowledgement 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 acknowledgement is made in the text. Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my thesis/dissertation, 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

The above candidate has carried out research for the Masters thesis/Dissertation under my supervision.

Signature of the Supervisor:                                 Date:

Name of the Supervisor: Dr. Amal Shehan Perera

# ABSTRACT

## A heterogeneous data ensemble approach for the classification of *Saccharomyces cerevisiae* proteins under 'mitochondrion organization'

Proteins are the real role players in keeping a cell healthy and well functioning. An important group of proteins is the subset of mitochondrial proteins that engage in the assembly, arrangement and disassembly of the mitochondrion. Several of them have been identified to cause human diseases. Hence, annotating proteins under the 'mitochondrion organization' Biology process is vital for identifying disease causative factors and for designing therapeutics. As manual annotation requires costly and laborious in vitro methods, in silico function prediction is preferred nowadays. Recent studies identify the importance of incorporating data from various biological aspects, to formulate a strong functional context for classification. In addition, many approaches from literature employ ensemble classifiers to attain a higher prediction accuracy. However, an insightful approach for accurate classification; biological data utilization; and biological data type significance determination; is still in need. This study presents an assessment of a heterogeneous data ensemble to classify *Saccharomyces cerevisiae* proteins under 'mitochondrion organization'. The ensemble consists of nine euclidean-distance based nearest neighbour models and three affinity-based neighbourhood models; it utilizes sequences, protein domains, peptide chain properties, gene expression, secondary structure and interactions. The base models were trained upon annotations from the Gene Ontology, as well as from a publicly available benchmark gold dataset. They show a substantial level of disagreement, implying their effectiveness in collective decision making. Six combination schemes were evaluated for fusing the base model outputs. A Genetic Algorithmically weighted ensemble gives the highest improvement to the best performing base classifier, by displaying an average area under the Receiver Operating Characteristic curve of 92.52%. Moreover, it is capable of determining the biological importance of each data type. Overall, the proposed heterogeneous data ensemble is capable of identifying eight disease related proteins and one disease related protein in a strong and moderate sense, respectively.

**Keywords**: yeast; proteins; mitochondrion; weighted ensemble; data heterogeneity; genetic algorithm; supervised learning

To my beloved parents, grandmother and brother

# ACKNOWLEDGEMENT

**I would like to express my heartiest and sincere gratitude,**

To my parents, for all their support, guidance, motivation and inspiration

To my advisor and supervisor Dr. Amal Shehan Perera, for his immense support, invaluable advice, continuous guidance and encouragement, through productive discussions and progress reviews, in making this research a success

To my Research Review Committee: Prof. Nalin Wickramarachchi and Dr. Dulani Meedeniya for their constructive feedback and encouragement

To Prof. T. L. Shamala Tirimanne from the University of Colombo, for offering me with her expertise in Biology through informative discussions, despite her busy schedule

To Dr. Surangika Ranathunga and Dr. Charith Chitraranjan, for those illuminating and motivating discussions despite their busy schedules

To Prof. Gihan Dias, for his constant advice and guidance

To Prof. Vajira H. W. Dissanayake and Mr. Nilaksha Neththikumara from the Human Genetics Unit, University of Colombo, for providing me with Training in Bioinformatics

To the Department of Computer Science and Engineering, the Senate Research Grant Committee, the Faculty of Graduate Studies and the staff in general at the University of Moratuwa, for supporting and facilitating my research with necessary resources throughout the course of study

# TABLE OF CONTENTS

University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| AGPS | Annotating Genes with Positive Samples |
| ANOVA | Analysis of Variance |
| AUC | Area Under the Curve |
| BioGRID | Biological General Repository for Interaction Datasets |
| BLAST | Basic Local Alignment Search Tool |
| CAFA | Critical Assessment of protein Function Annotation |
| CD | Czekanowski-Dice |
| CTD | Conjoint Triad Descriptor |
| Da | Dalton (atomic mass unit) |
| DF | Degrees of Freedom |
| DNA | Deoxyribonucleic Acid |
| FunCat | Functional Catalogue |
| GA | Genetic Algorithm |
| GO | Gene Ontology |
| GPCR | G Protein-Coupled Receptor |
| HER2 | Human Epidermal Growth Factor Receptor 2 |
| IEA | Inferred from Electronic Annotation |
| LDA | Latent Dirichlet Allocation |
| MIPS | Munich Information Center for Protein Sequences |
| NGS | Next Generation Sequencing |
| NLP | Natural Language Processing |
| NMR | Nucleic Magnetic Resonance |
| NN | Nearest Neighbour |
| mRNA | Messenger Ribonucleic Acid |

PAAC    Pseudo Amino Acid Composition

PCT     Predictive Clustering Tree

PDB     Protein Data Bank

PPI     Protein Protein Interactions

PR      Precision-Recall

QSOD    Quasi Sequence Order Descriptor

RNA     Ribonucleic Acid

ROC     Receiver Operating Characteristic

SGD     Saccharomyces Genome Database

SS      Secondary Structure

SVM     Support Vector Machine

TMC     Transductive Multi-label Classifier

TPR     True Path Rule

3D      Three dimensional