Ain Shams University
Faculty of Computer and Information Sciences
Information Systems Department



### Mining Structural Patterns for Automatic Protein Function Prediction

A thesis submitted in partial fulfillment of the requirements for the degree of PhD in Computer and Information Sciences

Tο

Information Systems Department,
Faculty of Computer and Information Sciences,
Ain Shams University

By

### **Huda Amin Maghawry Amin**

M.Sc. in Computer and Information Sciences,
Assistant Lecturer at Information Systems Departments
Faculty of Computer and Information Sciences, Ain Shams University

# Under Supervision of Prof. Dr. Mostafa Gadal-Haqq M. Mostafa

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

#### Prof. Dr. Mohamed Hashem Abdel Aziz

Information Systems Department Faculty of Computer and Information Sciences, Ain Shams University

#### Prof. Dr. Tarek Fouad Gharib

Information Systems Department Faculty of Computer and Information Sciences, Ain Shams University

### Acknowledgement

First and foremost, thanks to Allah, the most merciful, the most graceful, for His great help throughout this work.

I would like to express my deep and sincere gratitude to Prof. Dr. Mostafa Gadal-Haqq. His leadership, patience, understanding, encouragement, immense knowledge and personal guidance have provided a significant basis for this work. His detailed review and constructive comments have been of great value for me. I am gratefully thankful.

Also, I would like to express my deep and sincere gratitude to Prof. Dr. Tarek Gharib. I deeply appreciate his continuous supportive supervision. His advices, enthusiasm, patience and understanding were very motivational for me. I am gratefully thankful.

I wish also to express my warm and sincere thanks to Prof. Dr. Mohamed Hashem for his help and support.

My sincere and special gratitude from my deep heart is due to my great and beloved parents who supported me throughout my life with all their love and faith. They were always around at times I thought that it is impossible to continue, they helped me to keep things in perspective. Without their encouragement and understanding it would have been impossible for me to finish this work.

Last but not the least; I would like to thank my dear and beloved husband. His continuous and unconditional love, support, patience, sacrifice and understanding will remain my inspiration throughout my life. Also, my little and lovely son with his sweet smile, I owe him lots and lots of fun hours. Words would never say how grateful I am to both of them.

### **Abstract**

Protein function prediction is a challenging problem in bioinformatics. It has a great impact on the areas of diseases treatment and drug industry. Structure-based proteins representation plays an important role in proteins function prediction process. Three aspects of protein functions prediction have been considered: Predicting enzymes family and superfamily, classifying enzymes versus non-enzymes proteins and discriminating DNAbinding and non DNA-binding proteins. The thesis presents a modification to an existing protein representation approach which utilizes distance patterns between protein residues and a maximum cutoff. Also, the thesis presents a new protein structure representation for efficient protein function prediction. The new representation is based on three-dimensional patterns of protein residues. It utilizes atoms coordinates of protein residues, including the angles and distance patterns. This thesis also presents a study of the proposed protein representation and other protein-derived sequence, psychochemical and structure features to enhance the prediction of DNAbinding proteins and their classes.

## **Summary**

One of the challenging problems in bioinformatics is the prediction of protein function. Protein function is the main key that can be used to classify different proteins. The analysis of proteins and their functions is an important research area. Such analysis affects many applications like clarification of the living body mechanism, treatment of diseases and drug industry. Protein function can be inferred experimentally with very small throughput and high cost or computationally with very high throughput and lower cost. Many of protein sequences and structures are available but with no knowledge about their function. Therefore, methods for protein function prediction are highly and continuously required. Computational methods are based on protein sequences or structures. Protein functions are highly related to their structures. Therefore. structure-based proteins representation plays an important role in the prediction process.

This thesis presents a modification to an existing protein representation approach which utilizes distance patterns between protein residues and a maximum cutoff. The proposed modified representation considers the whole protein instead of using cutoff. Comparative analysis was done to evaluate the proposed representation method and the existing method. The aspect of protein function considered is based on enzyme activity. The results show that the proposed representation outperforms the existing representation with a prediction accuracy of 90.12% and 80.27% for superfamily and family level, respectively, with accuracy improvement of about 5% in average.

This thesis also presents a new structure-based protein representation for efficient protein function prediction. The new representation is based on three-dimensional patterns of protein residues. It utilizes atoms coordinates of protein residues, including the angles and distance patterns. The proposed representation uses protein structure only with no need to any sequence information. Besides, it does not need any prior alignment process. The aspects of protein functions considered using different datasets: Predicting enzymes family and superfamily and classifying enzymes versus non-enzymes proteins. The prediction accuracy of the proposed representation using various classification methods outperforms a recently introduced representation that is based only on the distance patterns. The results show that the proposed representation achieved prediction accuracy of 98.3% in predicting superfamily, 91% in predicting family and 79.25% in predicting enzyme proteins, with improvement of about 10% on average.

Finally, the thesis presents a study of different protein-derived sequence, psychochemical and structure features. The objective was to enhance the prediction of DNA-binding proteins and classes. This is achieved through finding efficient protein representations that able to predict whether a protein is DNA-binding protein and analyzing how well protein-derived representations predict each of DNA-binding protein classes. The protein features achieved accuracy improvement of about 7% on average for the prediction of DNA-binding proteins. The proposed representation when combined with other features achieved improvement in accuracy about 7% and 12% on average for the prediction of DNA-binding proteins and DNA-binding protein classes, respectively.

## **Table of Contents**

Acknowled	lgement	ii
Abstract		iii
Summary		iv
Table of Co	ntents	vi
List of Figu	res	ix
List of Tabl	es	xi
List of Abb	reviations	xiv
Chapter 1	Introduction	1
1.1 0	verview	1
1.2 M	lotivation	1
1.3 0	bjective	2
1.4 T	hesis Organization	2
Chapter 2	Background	4
2.1 Ir	ntroduction	4
2.2 B	ioinformatics	4
2.3 P	roteins	5
2.4 D	ata Mining	16
	1 Classification	
2.4.	2 Performance Measures	18

2.5 Sı	ımmary	20
Chapter 3	Related Works	21
	troduction	
3.2 Pi	rotein Representations	23
3.2.	1 Sequence-Based Protein Representations	24
3.2.	2 Structure-Based Protein Representations	26
3.3 M	ethodologies and Applications	28
3.3.	1 Sequence-Based Function Prediction Methods	29
3.3.	2 Structure-Based Function Prediction Methods	35
3.4 Pi	rotein Function Prediction Evaluation	38
3.5 Sı	ımmary	39
Chapter 4	Proposed Protein Representations	41
4.1 In	troduction	41
4.2 M	CSM Representation	42
4.3 A	Comparative Analysis of MCSM and CSM	48
	1 Dataset	
4.3.	2 Experiments Setup	48
4.3.	3 Results and Discussion	49
4.4 A	New Protein Structure Representation	57
4.4.	1 Protein Structure Matrix (PSM)	57
4.4.	2 Protein Structure Matrix with Cutoff (PSM-C)	59
4.5 E	valuation of PSM and PSM-C	64
4.5.	1 Datasets	64

4.5	3 Results and Discussion	65
4.6 S	ummary	73
Chapter 5	Enhanced Prediction of DNA-Bind	ing Proteins
and Classes	s 74	
5.1 Ir	ntroduction	74
5.2 M	lethodology	76
5.3 D	atasets and Experiments Setup	77
5.4 F	eature Selection	80
5.5 P	rediction of DNA-Binding Proteins	83
5.6 P	rediction of DNA-Binding Classes	98
5.7 S	ummary	107
Chapter 6	Conclusion and Future Work	108
6.1 M	ICSM Protein Representation	108
6.2 P	SM and PSM-C Protein Representations	109
	Enhanced Prediction of DNA-Binding	
6.4 F	uture Work	111
List of Pub	lications	113
References		114

# **List of Figures**

Figure 4.13 PSM Algorithm	60
Figure 4.14 The flowchart of PSM representation	61
Figure 4.15 PSM-C Algorithm	62
Figure 4.16 The flowchart of PSM-C representation	63
Figure 4.17 Recall comparison between CSM and PSM-C for	different
enzymes superfamilies	71
Figure 4.18 Precision comparison between CSM and PSM-C for	different
enzymes superfamilies	72
Figure 5.1 Sensitivity of each feature for P248 dataset	85
Figure 5.2 Sensitivity of each feature for P304 dataset	86
Figure 5.3 Sensitivity of each feature for P359 dataset	86

## **List of Tables**

Table 2.1 Conversion between single and three-letter amino acid codes and
RNA codon7
Table 2.2 Hydrophobic values of amino acids13
Table 3.1 Comparison of protein subcellular localization prediction methods
based on different sequence features using Reinhardt's dataset [72]33
Table 3.2 Comparison of protein subcellular localization prediction methods
using different sequence features on different datasets34
Table 4.1 Superfamily prediction accuracy using Random Forest and KNN
with non-zero minimum distance50
Table 4.2 Family prediction accuracy using Random Forest and KNN with
non-zero minimum distance50
Table 4.3 Superfamily prediction accuracy using Random Forest and KNN
with zero minimum distance54
Table 4.4 Family prediction accuracy using Random Forest and KNN with
zero minimum distance54
Table 4.5 Superfamily prediction accuracy using different classification
algorithms65
Table 4.6 Family prediction accuracy comparison using different
classification algorithms66
Table 4.7 Enzyme prediction accuracy comparison using different
classification algorithms67
Table 4.8 Comparison of superfamily prediction accuracy for PSM-C and CSM
with and without SVD68

Table 4.9 Superfamily classification comparison of CSM and PSM-C v	ısing
Naive Bayes	69
Table 4.10 Superfamily classification comparison of CSM and PSM-C $\alpha$	ısing
KNN	69
Table 4.11 Superfamily classification comparison of CSM and PSM-C $\iota$	ısing
Random Forest.	70
Table 5.1 Summary of proteins datasets	78
Table 5.2 Class distribution of DNA-binding proteins datasets	79
Table 5.3 Prediction accuracy using differnt feature selection methods u	ısing
benchmark datasets	81
Table 5.4 Number of correctly classified proteins using each feature for I	2248
datasetdataset	84
Table 5.5 Number of correctly classified proteins using each feature for I	2304
datasetdataset	84
Table 5.6 Number of correctly classified proteins using each feature for I	2359
datasetdataset	85
Table 5.7 Feature representations achieving maximum MCC for P248	87
Table 5.8 Feature representations achieving maximum MCC for P304	88
Table 5.9 Feature representations achieving maximum MCC for P359	89
Table 5.10 Prediction results comparison using P248 dataset	91
Table 5.11 Prediction results comparison using P304 dataset	92
Table 5.12 Prediction results comparison using P359 dataset	93
Table 5.13 DNA-binding proteins prediction results using PSM-C	94
Table 5.14 Number of correctly classified proteins using PSM-C	94
Table 5.15 Improvements achieved after adding PSM-C	95

Table 5.16 DNA-binding proteins prediction results with and without PSM-C
using Random Forest97
Table 5.17 DNA-binding proteins prediction results with and without PSM-C
using SVM97
Table 5.18 The sensitivity and the number of correctly classified proteins
using SVM98
Table 5.19 Correctly classified proteins using Random Forest for DB54
dataset99
Table 5.20 The sensitivity of each feature using Random Forest for DB54
dataset99
Table 5.21 Comparison of correctly classified proteins using DB54 dataset.
Table 5.22 Correctly classified proteins using PSM-C and Random Forest. 101
Table 5.23 Improvement indication after using PSM-C for RDB53 dataset.103
Table 5.24 Improvement indication after using PSM-C for R239 dataset 103
Table 5.25 Sensitivity comparison of each feature with and without PSM-C
for RDB53 dataset104
Table 5.26 Sensitivity comparison of each feature with and without PSM-C
for RDB239 dataset105
Table 5.27 DNA-binding protein classes prediction results with and without
PSM-C

### **List of Abbreviations**

2SAAC Two-segment Amino Acid Composition

AAC Amino Acid Composition

AACD Amino Acid Composition Distribution

ANN Artificial Neural Network

APGM Approximate Graph Mining

BLAST Basic Local Alignment Search Tool

CAFA Critical Assessment of protein Function Annotation

CDF Cumulative Distribution Function

CE Combinatorial Extension

CFS Correlation-based Featured Subset

CSM Cutoff Scanning Matrix

DALI Distance-matrix ALIgnment

DC Dipeptide Composition

DNA Deoxyribonucleic Acid

DWT Discrete Wavelet Transform

EC Enzyme Commission

ETA Evolutionary Trace Annotation

FAIR Finding All Internal Repeats

FASTA Fast-All

FN False Negative
FP False Positive

GAAP Gapped Amino Acid Pair Composition

GO Gene Ontology

GRAVY Grand Average of Hydropathicity Index

KNN K-Nearest Neighbor

MCC Mathew Correlation Coefficient

MCSM Modified Cutoff Scanning Matrix

MD Moment Descriptor

MIPS Munich Information Center for Protein Sequences

mRNA messenger Ribonucleic Acid

MSE Multi-Scale Energy

NMR Nuclear Magnetic Resonance

OIT Optimal and Information Theoretic

PCA Principle Components Analysis

PcAA Pair-coupled Amino Acid Composition

PDB Protein Data Bank

PDF Probability Distribution Function

PPI Protein-protein Interactions

PSAS Pairwise Sequence Alignments Scores

PseAAC Pseudo Amino Acid Composition

PSI-BLAST Position-Specific Iterative Basic Local Alignment Search Tool

PSM Protein Structure Matrix

PSM-C Protein Structure Matrix using Cutoff

RBF Radial Basis Function

ROC Receiver Operating Characteristics

SALSA Structurally Aligned Local Sites of Activity

SCOP Structural Classification of Proteins

SSAP Sequence Structure Alignment Program

SSMBS Sequentially Separated Motifs in Biological Sequences

SVD Singular Value Decomposition

SVM Support Vector Machine

TN True Negative

TOPS Topology-based Protein Structure

TP True Positive