

## بسم الله الرحمن الرحيم

 $\infty\infty\infty$ 

تم رفع هذه الرسالة بواسطة / حسام الدين محمد مغربي

بقسم التوثيق الإلكتروني بمركز الشبكات وتكنولوجيا المعلومات دون أدنى مسئولية عن محتوى هذه الرسالة.

AIN SHAMS UNIVERSITY

Since 1992

Propries 1992

ملاحظات: لا يوجد



# Ain Shams University Faculty of Computer and Information Sciences Information Systems Department Bioinformatics Program

## Artificial Intelligence Approach for Protein Sequence Analysis

A thesis submitted as partial fulfillment of the requirements for the degree of Master's in Information Systems, Bioinformatics program

By

#### Farida Alaaeldin Mostafa Mohamed

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

Under the Supervision of

#### Prof. Dr. Nagwa Badr

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

#### Prof. Dr. Rasha Ismail

Professor at Information Systems Department, Vice Dean for Graduate Studies and Research, Faculty of Computer and Information Sciences, Ain Shams University

#### **Dr. Yasmine Afify**

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

> 2022 Cairo

#### Acknowledgment

The assistance of my supervisors made this thesis achievable. I would want to thank Dr. Nagwa Badr, for her genuineness and support, which I will never forget. She is the icon of leadership and the ideal role model. Dr. Rasha Ismail has been an inspiration to me while I worked through this master's degree. This thesis would not have been possible without the help of Dr. Yasmine Afify, who guided me through the research process from the beginning and helped me gain a better grasp of the subject.

I am grateful for the tremendous possibilities my supervisors provided for me to progress professionally and for the extraordinary experiences they planned for me.

I am grateful for my parents' unwavering love and support, which keeps me driven and self-assured. My achievements and success are due to their belief in me. My gratitude goes out to my siblings, who remind me of the important things in life and are always encouraging in my experiences.

Finally, thank you to my friends for continually listening to me complain and chat things out, for cracking laughs when things grew too serious, and for making things better so that I could pursue a master's degree.

#### **Abstract**

Protein sequence analysis helps in the prediction of protein functions. The objective of this thesis is to propose new deep learning models that are capable of classifying proteins based on their features extracted in either 1D or 3D and investigate the impact of data variations using 3D features on the deep learning-based protein sequence classification.

Regarding the 1D features, different protein descriptors were used and decomposed into modified feature descriptors using Empirical Mode Decomposition that were not employed in protein studies. Uniquely, we introduced using Convolutional Neural Network to learn and classify protein diseases. A dataset of 1563 protein sequences was classified into 3 different disease classes: AIDS, Tumor suppressor, and Proto-oncogene.

Results showed a significant increase in the performance of the Convolutional Neural Network model using modified feature descriptor over Support Vector Machine using rbf kernel function by 23.3% in accuracy. CTDT modified feature descriptors improved the deep learning model results by 19.5%, 39.6%, 23.3%, 29.9%, 24.3%, and 31.2% in AUC, MCC, accuracy, F1- score, recall, and precision, evaluation metrices respectively.

Regarding the 3D features, uniquely five feature extraction groups were utilized to create 3D features with two sizes (7x7x7 and 9x9x9). Three datasets are employed in the assessment, which are different in their sorts, sizes, and balance state namely, Disease and two Phage Virion Proteins datasets.

Results showed that the 7x7x7 feature matrix has a positive correlation between its dimensions, which has positive impact on the results reaching 71% in PVP-Balanced and 86% in disease dataset. Using the sum of the first three Intrinsic Mode Function components had a better impact than using the first component improving accuracy to 86.6% for disease dataset. The dataset size had a significant positive impact on training the Convolutional Neural Network model reaching 84%.

#### **List of Publications**

- [1] Farida Alaaeldin Mostafa, Yasmine Mohamed Afify, Rasha Mohamed Ismail, Nagwa Lotfy Badr, "Protein Deep Learning Classification Using 3D Features," In 2021 Tenth International Conference on Intelligent Computing and Information Systems (ICICIS), 2021, pp. 462-466, DOI: 10.1109/ICICIS52592.2021.9694247
- [2] Farida Alaaeldin Mostafa, Yasmine Mohamed Afify, Rasha Mohamed Ismail, Nagwa Lotfy Badr, "Deep Learning Model for Protein Disease Classification" In Current Bioinformatics, 2022; 17(3), pp. 245-253, DOI: 10.2174/1574893616666211108094205
- [3] Farida A. Mostafa, Yasmine M. Afify, Rasha Ismail, Nagwa Badr. "Uncovering The Effects of Data Variation on Protein Sequence Classification Using Deep Learning" In International Journal of Intelligent Computing and Information Sciences, 2022; 22(2): 112-125. DOI: 10.21608/ijicis.2022.123177.1168

### **Table of Contents**

Acknowled	gment	II
Abstract		III
List of Publ	ications	IV
Table of Co	ontents	V
List of Figu	res	IX
List of Tabl	es	XI
List of Abb	reviations	XII
Chapter 1.	Introduction	2
1.1	Overview	2
1.2	Research Motivation	3
1.3	Research Objective	4
1.4	Research Contributions	5
1.5	Thesis Organization	5
Chapter 2.	Background and Related Work	7
2.1	Background	7
2.2	Related Work	10
	2.2.1 Analysis Using Graphical Representation	10
	2.2.2 Analysis Using Statistical Techniques	12
	2.2.3 Analysis Using Machine Learning Algorithms	13
	2.2.4 Analysis Using Deep Learning Algorithms	15
2.3	Summary	18
Chapter 3.	Protein Sequence Classification using Deep Learning	20
3.1	The Proposed Workflow Details	21
	3.1.1 Feature Extraction	21
	3.1.1.1 Amino Acid Composition (AAC)	23

	3.1.1.2	Grouped Di-Peptide Composition (GDPC)
	3.1.1.3	C/T/D Composition
	3.1.1.4	C/T/D Transition
	3.1.1.5	C/T/D Distribution
	3.1.1.6	Conjoint Triad
	3.1.1.7	Sequence-Order-Coupling Number (SOC Number)
	3.1.2 F	Feature Preprocessing
	3.1.3 N	Model Construction
	3.1.3.1	Convolutional Layer
	3.1.3.2	Dense Layers
	3.1.3.3	Drop out Layers
	3.1.3.4	Hyperparameters
	3.1.4 N	Model Training29
	3.1.5 F	Protein Classification
3.2	1D Syst	tem Architecture for Protein Sequence Analysis 32
	3.2.1	The 1D Feature Vector Extraction
	3.2.2	The 1D Feature Vector Preprocessing
	3.2.3	The 1D Model Constructing
	3.2.4	The 1D Model Training35
	3.2.5	The Protein Classification35
3.3	3D Syst	tem Architecture for Protein Sequence Analysis 36
	3.3.1	The 3D Feature Matrix Extraction
	3.3.2	The 3D Feature Matrix Preprocessing
	3.3.3	The 3D Model Construction
	3.3.4	The 3D Model Training
	3.3.5	The Protein Classification

3.4	Summary		
Chapter 4.	Experiments and Results		
4.1	Protein Datasets	42	
	4.1.1 Phage Virion Proteins datasets	42	
	4.1.2 Disease dataset	42	
4.2	Evaluation Metrics	43	
4.3	1D CNN Model Experiments	44	
	4.3.1 Experiment I	44	
	4.3.2 Experiment II	46	
	4.3.3 Experiment III	48	
	4.3.4 1D CNN Model Experiments Summary	51	
4.4	3D CNN Model Experiments	52	
	4.4.1 Experiment I	52	
	4.4.1.1 Results of PVP-Benchmark dataset	52	
	4.4.1.2 Results of PVP-Balanced dataset	53	
	4.4.1.3 Results of Disease dataset	53	
	4.4.2 Experiment II	55	
	4.4.3 Experiment III	58	
	4.4.4 Experiment IV	60	
	4.4.5 Experiment V	62	
	4.4.6 3D CNN Model Experiments Summary	62	
Chapter 5.	Conclusion and Future Work	65	
5.1	Conclusion	65	
5.2	Future Work	67	
References		69	

## **List of Figures**

Figure 2.1 Protein Structure Categories. (A) Primary Structure, (B) Secondary
Structure, and (C) Tertiary Structure
Figure 3.1 Proposed Workflow Diagram for Protein Classification 20
Figure 3.2 Proposed System Architecture for Protein Sequence Analysis.
(CONT.)
Figure 3.3 The Proposed 3D CNN Model Layers with Input and Output
Details
Figure 4.1 Accuracy comparison between SVM with four different kernel
functions
Figure 4.2 Performance Comparison of CNN, SVM-Rbf, and SVM Poly
Using Normal Feature Vectors with Several Evaluation Metrics: (A)
Precision, (B) MCC, (C) Accuracy, (D) F1-Score, (E) Recall, And (F) AUC.
47
Figure 4.3 Performance Comparison of CNN, SVM-Rbf, and SVM Poly
Using Modified Feature Vectors with Several Evaluation Metrics: (A)
Precision, (B) MCC, (C) Accuracy, (D) F1-Score, (E) Recall, And (F) AUC.
Figure 4.4 Evaluation Metrics on CNN Model Using Normal and Modified
Feature Vectors on: (A) AAC, (B) CTDT, (C) GDPC, and (D) CTDC 50
Figure 4.5 Performance Evaluation of the CNN Model PVP-Benchmark
Dataset
Figure 4.6 Performance Evaluation of the CNN Model on PVP-Balanced
Dataset

Figure 4.7 Performance Evaluation of the CNN Model on Disease Dataset.	
Figure 4.8 Loss Curves for Different Feature Matrix Sizes. (A) Loss Curve	
For 7x7x7, (B) Loss Curve For 9x9x9	

## **List of Tables**

Table 2.1 The 20 Amino Acids that Make Proteins
Table 3.1 The Five Groups of Amino Acid Feature Extraction Procedures. 22
Table 3.2 Detailed Description of the 8 Layers of the CNN Model Used for
Disease Classification. (~) Differs Based on The Input Size of the Feature
Descriptor
Table 4.1 Performance Evaluation on the Training Set on Three Datasets
Using Different Feature Matrix Sizes
Table 4.2 Performance Evaluation on the Independent Set on Three Datasets
Using Different Feature Matrix Sizes
Table 4.3 Performance Evaluation on the Training Set on Three Datasets
Using Different Components on IMFs on 7x7x7 Feature Matrix 59
Table 4.4 Performance Evaluation on the Independent Set on Three Datasets
Using Different Components on IMFs on 7x7x7 Feature Matrix 59
Table 4.5 Performance Evaluation on the Training Set on Three Datasets with
and without Normalization on 7x7x7 Feature Matrix
Table 4.6 Performance Evaluation on the Independent Set on Three Datasets
with and without Normalization on 7x7x7 Feature Matrix

#### **List of Abbreviations**

AAC Amino Acid Composition

ASMF Adaptive Signal Model Fingerprinting

AUC Area Under the Curve

CNN Convolutional Neural Network

CTDC C/T/D Composition

CTDT C/T/D Transition

CTriad Conjoint Triad

DFT Discrete Fourier Transform

DWT Discrete Wavelet Transform

ECG ElectroCardioGram

EMD Empirical Mode Decomposition

FN False Negative

FP False Positive

GAAC Grouped Amino Acid Composition

IMFs Intrinsic Mode Functions

MCC Matthews Correlation Coefficient

MSA Multiple Sequence Alignment

MSE Mean Square Error

NADH Nicotinamide Adenine Dinucleotide

ND5 Nicotinamide Adenine Dinucleotide dehydrogenase 5

ND6 Nicotinamide Adenine Dinucleotide dehydrogenase 6

Poly Polynomial

PRD Percent Root mean square Difference

PVP Phage Virion Proteins

PVPred-SCM Phage Virion Prediction-Scoring Card Method

RBF Radial Basis Function

ReLU Rectified Linear activation function

RQA Recurrence Quantification Analysis

SCM Scoring Card Method

SNRimp improved Signal to Noise Ratio

SVM Support Vector Machine

SVM-PCD Support Vector Machine Physicochemical Distributions

SVM-RQA Support Vector Machine Recurrence Quantification

TN True Negative

TP True Positive

## **Chapter 1**

# Introduction