

شبكة المعلومات الجامعية التوثيق الإلكتروني والميكروفيلو

# بسم الله الرحمن الرحيم





HANAA ALY



شبكة المعلومات الجامعية التوثيق الإلكتروني والميكرونيله



شبكة المعلومات الجامعية التوثيق الالكتروني والميكروفيلم



HANAA ALY



شبكة المعلومات الجامعية التوثيق الإلكترونى والميكروفيلم

## جامعة عين شمس التوثيق الإلكتروني والميكروفيلم قسم

نقسم بالله العظيم أن المادة التي تم توثيقها وتسجيلها على هذه الأقراص المدمجة قد أعدت دون أية تغيرات



يجب أن

تحفظ هذه الأقراص المدمجة بعيدا عن الغبار



HANAA ALY



#### **Ain Shams University**

## **Faculty of Computer and Information Sciences Information Systems Department**

## NAMED ENTITY RECOGNITION FROM BIOMEDICAL TEXT

Thesis submitted as a partial fulfillment of the requirements for the degree of Master of Science in Computer and Information Sciences

By

#### **Lobna Ahmed Mady**

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

#### **Under Supervision of**

#### Prof. Dr. Nagwa Lotfy Badr

Dean of Faculty of Computer and Information Sciences,
Ain Shams University

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

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

#### Acknowledgment

I would like to express my sincere gratitude and appreciation to my research main supervisor, Prof. Nagwa Badr. She provided me with valuable guidance and advice throughout this research work. It was really a great opportunity for me to gain from her deep knowledge.

I would like to express my appreciation to my research co-supervisor, Dr. Yasmine Afify, for her valuable comments, support, and encouragement throughout this research. It has been a privilege to work under her supervision and to learn from her experience. This work could not be done without her support.

I would like to thank my backbone and great father, Ahmed Mady, for all the guidance and advice through my whole life. Also, I would like to thank my wonderful mother, Nabwia Sabir, for providing the greatest support and help ever for me as well as for always trying to make my life easy, excellent, and happy.

I would also thank my husband and my other half, Dr. Ahmed Elansary, for being supportive during my master's journey and for always being ready to offer any kind of help.

Also, I want to thank my awesome sister, Soha Mady, for being that kind and caring for me and for all the times she tried to cheer me up and continue my work. More importantly, I have to thank my old brother, Mohamed Mady, for providing help and for all the times he pushed me to continue during hard times. Many thanks also should go to Dr. Abdellatif for his valuable assistance that he provided. Last but not the least, a very special thanks goes to the best friends ever Mariam and Nouran for always being beside me and providing support during whole course of my research.

#### **Abstract**

Named entity recognition is commonly considered an essential task in natural language processing. It represents a major phase in information extraction methodology to discover the named entities referenced in unstructured text and classify them into predefined class labels. Identifying biomedical entities has been recognized as a challenging task in named entity recognition.

In this thesis, the applicability of using structured support vector machine to classify flat and nested biomedical entities combined with the feature selection techniques to enhance the performance of biomedical named entity recognition has been thoroughly investigated. The proposed approach used a combination of various types of features to explore the classification performance in a combination of structured support vector machine as a machine learning technique. These features include linguistic, morphological, orthographical, context, and word representation features. The experimental results showed that the performance of the proposed approach surpassed that produced from other benchmark approaches in extracting the biomedical entities such as genes, proteins, cell lines, cell types, DNAs and RNAs.

Derived by these promising results, we were motivated to explore the effect of different types of features on structured support vector machine performance in extracting the biomedical entities. This was achieved by applying two types of filter-based feature selection techniques. The Chisquared and ReliefF feature selection techniques were utilized to reduce the feature vector and evaluate the importance of each feature. The reduced

feature vector was then taken as an input to structured support vector machine to extract the biomedical entities. Experimental results revealed that the overall performance of biomedical named entity recognition increased when using the adopted feature selection techniques and removed the features with the least effect on the classes.

Most studies in the literature ignored the nested entities and focused on the extraction of flat named entities only. However, the nested entities are commonly used in real world biomedical applications due to their ability to represent semantic meaning of the named entity. Therefore, the proposed approach was evaluated on the extraction of the nested biomedical entities and showed very promising results compared to those obtained from the benchmark approaches.

Comprehensive experiments were conducted on three popular datasets based on five evaluation metrics namely: recall, precision, F1-measure, Geometric Mean, and Matthews correlation coefficient. Experimental results revealed that the structured support vector machine achieved better performance compared to different approaches in the literature for extraction of both flat and nested entities.

### **Table of Contents**

Acknowled	gment	2
Abstract		3
Table of Co	ontents	5
List of Tabl	les	8
List of Figu	res	9
List of Abb	reviations	10
List of Publ	lications	12
Chapter 1	Introduction	13
1.1	. Overview	13
1.2	2. Problem Definition	14
1.3	3. Research Objective	14
1.4	4. Contribution	14
1.5	5. Thesis Organization	15
Chapter 2	Background	16
2.1	Named Entities (NEs)	16
2.2	2. Named Entity Recognition (NER)	16
2.3	Biomedical Data	17
2.4	4. Biomedical Named Entity Recognition (BNER)	17
2.5	5. Feature Engineering	19
2.6	6. Model Selection	19
2.7	7. Machine Learning Techniques (ML)	20

2.7.1 Conditional Random Field (CRF)	20
2.7.2 Support Vector Machine (SVM)	21
2.7.3 Artificial Neural Networks (ANNs)	22
2.8. Feature Selection	23
2.8.1. Wrapper Approaches	24
2.8.2. Filter Approaches	25
2.8.2.1. Chi-squared feature selection technique	26
2.8.2.2. ReliefF feature selection technique	27
Chapter 3 Related Work	29
3.1. BNER Approaches	29
3.2. BNER Approaches Using Feature Selection Technique	34
3.3. Nested BNER Approaches	37
Chapter 4 Proposed Approach	40
4.1. Feature Engineering	40
4.1.1. Linguistic Features	40
4.1.2. Morphological Features	40
4.1.3. Orthographical Features	41
4.1.4. Context Features	41
4.1.5. WR Features	41
4.1.5.1. WE Feature	41
4.1.5.2. Clustering based Feature	43
4.2. Model Selection	49

4.3. Nested Biomedical Named Entity Recognition (Nest	ted
BNER)	50
Chapter 5 Experimental Results and Discussion	52
5.1. Datasets	52
5.1.1. GeneTag Dataset	52
5.1.2. JNLPBA Dataset	52
5.1.3. Genia Dataset	53
5.2. Evaluation Metrics	54
5.3. Experiments	56
5.3.1. Experiment I	56
5.3.2. Experiment II	57
5.3.3. Experiment III	60
5.3.4. Experiment IV	63
Chapter 6 Conclusion	66
Chapter 7 Future Work	68
References	69

### **List of Tables**

Table 1 Examples of WE Vector of JNLPBA Tokens
Table 2 Description of Features Used in the Proposed Approach 44
Table 3 Frequencies for NEs in JNLPBA Dataset
Table 4 Number of Tokens in Each Class Label in Genia Dataset 54
Table 5 Performance Comparisson for GeneTag Dataset
Table 6 Performance Comparison for JNLPBA Dataset
Table 7 Evaluation Metrics for Different Entity Types in JNLPBA and
GENETAG Datasets Using SSVM
Table 8 Comparison between overall Performance Comparison for GENIA
Dataset Before and After Applying the Feature Selection Techniques 61
Table 9 Performance of Each Class Label without Applying Feature
Selection
Table 10 Performance of Each Class Label After Applying Chi-Squared
Feature Selection
Table 11 Performance of Each Class Label After Applying Relieff Feature
Selection63
Table 12 The Overall Recall, Precision and F1-Measure for Each Class
Label of nested entities in GENIA Using SSVM
Table 13 Performance Comparison of the Proposed Approach Against
Benchmark Approaches

### **List of Figures**

Figure 1 Named Entity Recognition Example	17
Figure 2 Biomedical Named Entity Recognition Example	18
Figure 3 Illustration of SVM Model	21
Figure 4 Feed Forward ANN	23
Figure 5 Backpropagation ANN	23
Figure 6 Wrapper Feature Selection	25
Figure 7 Filter Feature Selection	26
Figure 8 A sub tree generated by Brown Clustering for a group of	tokens
extracted from JNLPBA corpus	44
Figure 9 Example of Nested Entities	51

#### List of Abbreviations

AI Artificial Intelligence

ANN Artificial Neural Network

BA Boundary Assembly

Bidirectional Contextual Clues Named Entity

BCC-NER Recognition

Bidirectional Encoder Representations from

BERT Transformers

BiLSTM Bidirectional Long Short-Term Memory
BNER Biomedical Named Entity Recognition

CNNs Convolutional Neural Networks

Convolutional Neural Network and Long Short-Term

CNN-LSTM Memory

CRF Conditional Random Field DNA Deoxyribonucleic Acid

FLSTM Forward Long Short-Term Memory

FN False Negative
FP False Positive
G-mean Geometric Mean

GPRO Gene and Gene Product
HMMs Hidden Markov Models
IE Information Extraction

Knowledge enhanced Biomedical pretrained Language

KeBioLM Model

LSTM Long Short-Term Memory Networks

MCC Matthews correlation coefficient
MEMMs Maximum Entropy Markov Models

ML Machine Learning

MultiCNN Multiple Convolutional Neural Networks

NER Named Entity Recognition

NEs Named Entities

Nested BNER Nested Biomedical Named Entity Recognition

NLP Natural Language Processing

PCA Principal Component Analysis

POS Part of Speech

PSO Particle Swarm Optimization

RNA Ribonucleic Acid

SPBA Statistical Principle Based Approach
SSVM Structured Support Vector Machine

SVMs Support Vector Machines

TCN Temporal Convolutional Network

Temporal Convolutional Network with a Conditional

TCN-CRF Random Field True Positive

UMLS Unified Medical Language System

WE Word Embedding
WR Word Representation

#### **List of Publications**

Lobna A. Mady, Yasmine M. Afify, Nagwa L. Badr, "Biomedical Named Entity Recognition Using Structured Support Vector Machine," International Journal of Computers and Their Applications, Vol. 28, No. 4, pp. 222-229, 2021.

Lobna A. Mady, Yasmine M. Afify, Nagwa L. Badr, "Enhancing Performance of Biomedical Named Entity Recognition," Accepted in International Conference on Intelligent Computing and Information Systems, ICICIS 2021, IEEE, 2021.

Lobna A. Mady, Yasmine M. Afify, Nagwa L. Badr, "Nested Biomedical Named Entity Recognition," Accepted in International Journal of Intelligent Computing, and Information Sciences, 2021.