



Prediction of Nitrotyrosine Sites Based On Composition and Position Based Features

By

Ahmad Waseem (15026050008)

A Dissertation Submitted For the Degree of Master in Computer Science in the Department of Computer Science, School of Systems and Technology, University of Management and Technology.

May, 2018

Supervised by:

Dr. Yaseer Daniyaal Khan

FINAL APPROVAL

It is certified that the research work presented in this thesis entitled “Prediction of Nitrotyrosine Sites Based on Composition and Position Based Features” was conducted by Ahmad Waseem Ghauri under the supervision of Dr. Yaser Daanial Khan at University of Management and Technology, Lahore, Pakistan in May 2018 to fulfil the requirements of the degree in M.S/M.Phil in Computer Science.

- **Supervisor**

Dr. Yaser Daanial Khan

Associate Professor, Chairperson Department of Computer Science, School of Systems and Technology, Department of Computer Science,

University of Management and Technology

- **Director Graduate Studies**

Dr. Shoaib Farooq

Associate Professor, Director Graduate Studies, School of Systems and Technology, Department of Computer Science,

University of Management and Technology

DECLARATION

I hereby **Ahmad Waseem Ghauri** with Student ID **15026050008** declare in this declaration that it's my own original work. Where collaborations with other researchers are involved, or materials generated by other researchers are included, the parties and/or materials are acknowledged or are explicatory referenced as appropriate.

This work is being submitted to the University of Management and Technology in partial fulfilment of the requirements of the degree of Master of Science (M.S) in Computer Science.

This thesis has not been submitted to any other University or institution for any degree or examination.

Date

Student Signature

Date

Supervisor Signature

ABSTRACT

Closely related to causes of various diseases such as rheumatoid arthritis, septic shock, and coeliac disease; tyrosine nitration is considered as one of the most important post-translational modification in proteins. Inside a cell, such modifications occur accurately by the action of sophisticated cellular machinery. This task is accomplished by specific enzymes present in endoplasmic reticulum. The identification of potential tyrosine residues in a protein primary sequence which can be nitrated is a challenging task. To counter the prevailing, laborious and time-consuming experimental approaches, here we introduce a novel computational model. Based on experimentally verified tyrosine nitration sites, they are transformed to their feature vectors. An adaptive training algorithm is then used to train a back propagation neural network for prediction purposes. To objectively measure the accuracy of the proposed model, rigorous verification and validation tests are carried out which led to a promising accuracy of 88%, a sensitivity of 85% and a specificity of 89.18% and Mathew correlation coefficient of 0.627. We believe that this computational model may provide foundation for further investigation and can be used deal with the other PTM sites in proteins.

ACKNOWLEDGEMENTS

Firstly, I'd like to thank my supervisor, Dr. Yaser Daanial Khan for his for his patience, motivation, enthusiasm, and on-going support throughout my M.S study and research.

I'd would also like to thank Dr. Nauman Rasool for his encouragement and insightful comments.

Last but not the least, my parents deserve credit for always being there for me and providing their wise counsel.

TABLE OF CONTENTS

| | |
|---|------------|
| ABSTRACT | i |
| ACKNOWLEDGEMENTS | ii |
| TABLE OF CONTENTS | iii |
| LIST OF TABLES | vi |
| LIST OF FIGURES | vii |
| | |
| 1 Introduction | 1 |
| 1.1 Problem Definition | 1 |
| 1.2 Research Objectives..... | 1 |
| 1.3 Scope..... | 2 |
| 1.4 Significance and Impact on Society..... | 3 |
| 2 Background..... | 4 |
| 2.1 Protein..... | 4 |
| 2.1.1 Structure Hierarchy..... | 4 |
| 2.1.2 Protein Data Representation | 6 |
| 2.1.3 Text Strings..... | 7 |
| 2.2 Post Translational Modification..... | 8 |
| 2.2.1 An Introduction..... | 8 |

| | | |
|-------|--|----|
| 2.2.2 | Most Common PTMs | 9 |
| 2.2.3 | Influence on Diseases and Health..... | 9 |
| 2.3 | Nitration of Tyrosine | 10 |
| 2.3.1 | Origins | 10 |
| 2.3.2 | Significance in Biology | 11 |
| 3 | Literature Review | 13 |
| 4 | Materials and Methods | 15 |
| 4.1 | Dataset Collection..... | 15 |
| 4.1.1 | Program Code | 16 |
| 4.2 | Feature Vector Construction..... | 17 |
| 4.3 | Statistical Moments of Primary Structure..... | 18 |
| 4.4 | Position Relative Incidence Matrix..... | 22 |
| 4.5 | Frequency Matrix..... | 22 |
| 4.6 | Accumulative Absolute Position Incidence Vector (AAPIV) | 23 |
| 4.7 | Reverse Accumulative Absolute Position Incidence Vector (RAAPIV).. | 24 |
| 4.8 | Prediction Algorithm | 24 |
| 4.9 | Gradient Descent and Adaptive Learning..... | 25 |

| | | |
|-------|---|----|
| 5 | Results and Discussion | 29 |
| 5.1 | Results..... | 31 |
| 5.1.1 | Self-Consistency Test | 31 |
| 5.1.2 | Demonstration on Independent Set..... | 32 |
| 5.1.3 | Jackknife Cross Validation | 34 |
| 5.1.4 | K Fold Cross Validation | 35 |
| 5.2 | Discussion..... | 43 |
| 5.2.1 | Comparison with Existing Predictors | 43 |
| 6 | Conclusion and Implications | 45 |
| 7 | Future Research | 46 |
| 8 | References | 47 |

LIST OF TABLES

| | |
|---|----|
| Table 5.1. 10 Fold Cross Validation Results on our benchmark dataset..... | 42 |
| Table 5.2. 10-fold Cross Validation Results on the benchmark dataset used in [38]..... | 42 |
| Table 5.3. Test results derived by the 5-fold cross-validation on the benchmark dataset in [38]. | 43 |
| Table 5.4. Comparison of our model with the existing predictors | 44 |

LIST OF FIGURES

| | |
|--|----|
| Figure 2.1. Sample FASTA File | 7 |
| Figure 4.1. Y' (Tyrosine) as the centre..... | 16 |
| Figure 4.2. Neural Network | 25 |
| Figure 4.3. Prediction Process Model | 27 |
| Figure 4.4. Positive Peptide Frequency Plot | 27 |
| Figure 4.5. Negative Peptide Frequency Plot..... | 28 |
| Figure 5.1. Confusion Matrix - Self Consistency Test..... | 31 |
| Figure 5.2. Regression Plot - Self Consistency Test..... | 31 |
| Figure 5.3. Regression Plot - Self Consistency Test..... | 32 |
| Figure 5.4. ROC Curve of Independent Set | 33 |
| Figure 5.5. 1st Fold Test Confusion Matrix..... | 37 |
| Figure 5.6. 2nd Fold Confusion Matrix | 37 |
| Figure 5.7. 3rd Fold Confusion Matrix | 38 |
| Figure 5.8. 4th Fold Confusion Matrix | 38 |
| Figure 5.9. 5th Fold Confusion Matrix | 39 |
| Figure 5.10. 6th Fold Confusion Matrix | 39 |
| Figure 5.11. 7th Fold Confusion Matrix | 40 |
| Figure 5.12. 8th Fold Confusion Matrix | 40 |
| Figure 5.13. 9th Fold Confusion Matrix | 41 |
| Figure 5.14. 10th Fold Confusion Matrix | 41 |

1 Introduction

1.1 Problem Definition

Given the post genomic age we currently live in, the protein sequences generated in this era has been exponentially increased [1]. Nitrotyrosine is considered a marker for cell damage and inflammation and is linked to wide range of human pathological conditions [2] and is detected in a large number of diseases, such as lung cancer, cardiovascular disease, asthma, Alzheimer's disease, rheumatoid arthritis, septic shock and coeliac disease [3]. The determination of nitrotyrosine residues in a protein is not an easy and inexpensive task, but time consuming involving advanced technologies. The lack of experimentally verified susceptible sites has made it difficult to examine global biophysical and evolutionary trends of nitrotyrosine in a protein [4].

1.2 Research Objectives

To counter the prevailing, laborious and time-consuming experimental approaches for the detection of tyrosine nitration sites, our research objective is to introduce a novel computational model which is based on experimentally verified tyrosine nitration sites. Our aim is to predict nitrotyrosine sites which can provide us insights of its impact at the proteome level. In order to predict tyrosine residues which are prone to nitrosylation, many computational models have been proposed with varying sensitivities and accuracies. The prediction of nitrotyrosine sites in proteins is under active research by the community and considered basis for detection of pathological conditions and drug development [5].