

# **Email Classification Tool to Detect Phishing Using Hybrid Features**

H.V. Mahesh  
169321L

Faculty of Information Technology

University of Moratuwa

March 2019

# **Email Classification Tool to Detect Phishing Using Hybrid Features**

H.V. Mahesh  
169321L

Dissertation submitted to the Faculty of Information Technology, University of Moratuwa, Sri Lanka for the partial fulfillment of Degree of Master of Science in Information Technology.

**February 2019**

# Declaration

We declare that this thesis is our own work and has not been submitted in any form for another degree or diploma at any university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged in the text and a list of references is given.

Name of Student (s)

Signature of Student (s)

H. V. Mahesh

.....

Date .....

Supervised by

Name of Supervisor

Signature of Supervisor

S.C. Premarathne

.....

Date: .....

## **Acknowledgements**

First and foremost, I would like to express my sincere gratitude towards my supervisor, Mr. Saminda Premarathne, Senior Lecturer, Faculty of Information Technology, University of Moratuwa for his guidance, supervision, advices and sparing valuable time thorough the research project.

A special thanks goes to Mr. Indika Gunawardhana, CIO, LAUGFS Holdings Ltd for the valuable support and guidance given for the research. I would also like to thank all batch mates of my MSc IT program who gave their valuable feedbacks to improve the results of the research. Further I express my warm thanks to my family & wife for the support and encouragement they have given during this time.

## **Abstract**

Phishing is a fraudulent attempt of trying to gather personal sensitive information such as user ID and passwords, credit card and bank account details through network. Social messaging and websites are used as medium to trigger attacks in addition to the use of emails, which is the most common and leading method currently used to perform phishing attacks. In an attack, the attacker is sending an email with a URL of the phishing website camouflaged as a legitimate source.

Nowadays phishing has become more complicated and critical problem to many organizations. The phishers can bypass the filters and rules set by anti-phishing procedures and techniques. This research build a web based phishing email detection tool using data mining classification model.

To build an efficient classification model, varieties of extracted email features have been used. These selected features can be categorized according to email header, email body, URL and Web Page Content of URL. In this model, classification accuracy will be enhanced by using these hybrid features.

This model will be used to implement the web-based tool to detect phishing emails with more accuracy even without opening the emails. This can be used as preventive and proactive technique for phishing detection.

# Contents

|                                                   |      |
|---------------------------------------------------|------|
| DECLARATION .....                                 | I    |
| ACKNOWLEDGEMENTS .....                            | II   |
| ABSTRACT .....                                    | III  |
| CONTENTS .....                                    | IV   |
| LIST OF FIGURES .....                             | VII  |
| LIST OF TABLES .....                              | VIII |
| <b>CHAPTER 1 - INTRODUCTION</b> .....             | 1    |
| 1.1 PROLEGOMENA.....                              | 1    |
| 1.2 BACKGROUND AND MOTIVATION .....               | 1    |
| 1.3 PROBLEM STATEMENT .....                       | 2    |
| 1.4 AIMS AND OBJECTIVES .....                     | 2    |
| 1.4.1 Aim.....                                    | 2    |
| 1.4.2 Objectives.....                             | 3    |
| 1.5 PROPOSED SOLUTION.....                        | 3    |
| 1.6 STRUCTURE OF THE THESIS.....                  | 3    |
| <b>CHAPTER 2 - LITERATURE REVIEW</b> .....        | 4    |
| 2.1 INTRODUCTION.....                             | 4    |
| 2.2 RELATED WORK OF PHISHING CLASSIFICATION ..... | 4    |
| 2.3 SUMMARY OF RELATED STUDIES .....              | 9    |
| 2.4 LIST OF FEATURES.....                         | 10   |
| 2.4.1 Email Header Based.....                     | 11   |
| 2.4.2 Email Body Based.....                       | 11   |
| 2.4.3 URL Based.....                              | 12   |
| 2.4.4 URL Web Page Content Based .....            | 13   |
| 2.5 SUMMARY .....                                 | 14   |

|                                                                |           |
|----------------------------------------------------------------|-----------|
| <b>CHAPTER 3 - TECHNOLOGIES AND TOOLS USED FOR PHISHING</b>    |           |
| <b>CLASSIFICATION .....</b>                                    | <b>15</b> |
| 3.1 INTRODUCTION.....                                          | 15        |
| 3.2 DATA MINING TECHNIQUES .....                               | 15        |
| 3.3 NAIVE BAYES.....                                           | 15        |
| 3.4 K-NEAREST NEIGHBORS .....                                  | 16        |
| 3.5 DECISION TREE .....                                        | 16        |
| 3.6 RAPID MINER STUDIO.....                                    | 16        |
| 3.7 .NET FRAMEWORK .....                                       | 17        |
| 3.8 MICROSOFT VISUAL STUDIO .....                              | 17        |
| 3.9 EAGETMAIL.....                                             | 17        |
| 3.10 HTMLAGILITYPACK.....                                      | 17        |
| 3.11 TALEX.SEOSTATS .....                                      | 18        |
| 3.12 PHISHTANK API.....                                        | 18        |
| 3.13 MICROSOFT SQL SERVER.....                                 | 18        |
| 3.14 SUMMARY .....                                             | 18        |
| <br>                                                           |           |
| <b>CHAPTER 4 - A NOVEL APPROACH OF CLASSIFICATION PHISHING</b> |           |
| <b>EMAIL USING HYBRID FEATURES .....</b>                       | <b>19</b> |
| 4.1 INTRODUCTION.....                                          | 19        |
| 4.2 HYPOTHESIS .....                                           | 19        |
| 4.3 INPUT .....                                                | 19        |
| 4.4 OUTPUT .....                                               | 20        |
| 4.5 PROCESS AND FEATURES .....                                 | 20        |
| 4.6 USERS .....                                                | 20        |
| 4.7 SUMMARY .....                                              | 20        |
| <br>                                                           |           |
| <b>CHAPTER 5 - DESIGN OF THE CLASSIFICATION TOOL .....</b>     | <b>21</b> |
| 5.1 INTRODUCTION.....                                          | 21        |
| 5.2 HIGH LEVEL ARCHITECTURE OF SYSTEM.....                     | 21        |
| 5.3 DATA COLLECTION AND PREPROCESSING .....                    | 22        |
| 5.4 DESIGN OF CLASSIFICATION TOOL.....                         | 22        |

|                                                                |           |
|----------------------------------------------------------------|-----------|
| 5.5 BACK END DATABASE.....                                     | 23        |
| 5.6 SUMMARY .....                                              | 23        |
| <b>CHAPTER 6 - IMPLEMENTATION OF CLASSIFICATION TOOL .....</b> | <b>24</b> |
| 6.1 INTRODUCTION.....                                          | 24        |
| 6.2 CORE SERVICE.....                                          | 24        |
| 6.3 DATA COLLECTION BY DATA EXTRACTION TOOL.....               | 24        |
| 6.4 CLASSIFICATION MODEL .....                                 | 25        |
| 6.5 CLASSIFICATION TOOL.....                                   | 25        |
| 6.6 SUMMARY .....                                              | 26        |
| <b>CHAPTER 7 - EVALUATION.....</b>                             | <b>27</b> |
| 7.1 INTRODUCTION.....                                          | 27        |
| 7.2 EVALUATION OF CLASSIFICATION TECHNIQUES.....               | 27        |
| 7.3 SUMMARY .....                                              | 29        |
| <b>CHAPTER 8 - CONCLUSION AND FURTHER WORK .....</b>           | <b>30</b> |
| 8.1 INTRODUCTION.....                                          | 30        |
| 8.2 LIMITATIONS .....                                          | 30        |
| 8.3 FUTURE DEVELOPMENTS.....                                   | 30        |
| 8.4 SUMMARY .....                                              | 31        |
| <b>REFERENCES.....</b>                                         | <b>32</b> |
| <b>APPENDIX A - SAMPLE .EML FILE.....</b>                      | <b>34</b> |
| <b>APPENDIX B - MODEL EVALUATION SUMMARY .....</b>             | <b>35</b> |
| <b>APPENDIX C - DECISION TREE RULES.....</b>                   | <b>36</b> |
| <b>APPENDIX D – CODE SNIPPET.....</b>                          | <b>37</b> |
| <b>APPENDIX E – CLASSIFICATION TOOL UI.....</b>                | <b>39</b> |