

A COMPARISON AND EFFECTIVE PENETRATION TESTING APPROACHES WITH NMPREDICTOR BASED ON MACHINE LEARNING



By

MUHAMMAD NOMAN KHALID

A thesis

Presented to the Bahria University, Karachi Campus

In partial fulfillment of the requirement

For the degree of

MS (Computer Science)

SPRING, 2018

**DEPARTMENT OF COMPUTER SCIENCES
BAHRIA UNIVERSITY**

A COMPARISON AND EFFECTIVE PENETRATION TESTING APPROACHES WITH NMPREDICTOR BASED ON MACHINE LEARNING

By
Muhammad Noman Khalid

A thesis report submitted in partial fulfillment of the requirements for the degree of
MS. (Computer Science)

Supervisor: Dr. Humera Farooq

Nationality: Pakistani

Bahria University Karachi Campus
13, National Stadium Road
Pakistan



Thesis Completion Certificate

Student's Name: M. NOMAN KHALID Registration No. 24120

Programme of Study: MS (CS)

Thesis Title: A COMPARISON ON PENETRATION TESTING APPROACHES AND A PROPOSED FRAMEWORK FOR web Vulnerabilities BASED ON MACHINE LEARNING

It is to certify that the above student's thesis has been completed to my satisfaction and, to my belief, its standard is appropriate for submission for Evaluation. I have also conducted plagiarism test of this thesis using HEC prescribed software and found similarity index at 9% that is within the permissible limit set by the HEC for the MS/MPhil degree thesis. I have also found the thesis in a format recognized by the BU for the MS/MPhil thesis.

Principal Supervisor's Signature: _____

Date: 19/3/18 Name: Dr. Humera Farooq

Acknowledgement

First of all I would like to thank Almighty Allah who give me strength to complete my thesis. Then this study could not have been carried out without the cooperation, constructive advice and inspiration of my supervisor **Dr Humera Farooq**. I am equally grateful to **Sir Iqbal** who not only supported my idea but also gave me wise advises regarding my research. I would also like to thanks my colleague and friends who helped a lot to achieve goal of this research. Last but not the least I am thankful to my parents because their prayers and support are always with me.

Table of Contents

Acknowledgement	2
Table of Contents	3
List of Figures	6
List of Tables	7
List of Abbreviations	8
ABSTRACT	9
CHAPTER 1	10
INTRODUCTION	10
1.1 Background	10
1.2 Thesis Motivation	11
1.3 Problem Statement	12
1.4 Research objective	13
1.5 Contribution	13
1.5.1 Study on statistical classification based on existing method and to investigate the optimal method:	13
1.5.2 Proposed method (NMPREDICTOR):	14
1.6 Scope of Research	14
1.7 Thesis Organization	15
Chapter 2	16
Literature Review	16
2.1 Classification of web vulnerabilities	16
2.1.1 Improper input validation	17
2.1.1.1 Query manipulation	18
2.1.1.2 Client-side injection	18
2.1.1.3 File and path injection vulnerability	18
2.1.2 Improper authentication and authorization (Logic Flaw)	19
2.1.3 Improper session management	19
2.1.3.1 Session management	20
2.2 Detection of web vulnerabilities	20
2.3 Analysis Method to Detect Web Vulnerabilities	20
2.3.1 White Box Testing	20
2.3.2 Black box Testing	21

2.4	Static, Dynamic and Hybrid Analysis	21
2.5	Machine Learning Technique.....	23
2.6	Summary	31
CHAPTER 3		33
RESEARCH METHODOLOGY AND IMPLEMENTATION		33
3.1	Chapter Overview	33
3.2	Machine Learning and Web Vulnerability.....	33
3.3	Selected Papers for Comparative Analysis	34
3.3.1	Paper 1 P1: Usage of static analysis and data mining for the purpose of removal of Web application vulnerabilities (Medeiros et al., 2016).....	35
3.3.2	Paper 2 P2: Equipping of WAP with WEPONS for the detection of Vulnerabilities (Medeiros, Neves and Correia 2016).....	36
3.3.3	Paper 3 P3 Prediction of Vulnerable Components by Metrics vs. Text Mining (Walden, Stuckman and scandariato (2014).....	37
3.3.4	Paper 4 P4 an Empirical Investigation of Security Vulnerabilities within Web Applications (Abunadi and Alenezi, 2016).....	38
3.3.5	Paper 5 P5: For the prediction of Vulnerable File combination of Text Features and Software Metrics (YUN et al., 2015).....	38
3.3.6	Paper 6 P6: Prediction of Cross-Site Scripting (XSS) Security Vulnerabilities in Web Applications (Gupta et al., 2015).....	39
3.4	Proposed Method (NMPREDICTOR):.....	40
3.5	Evaluation Methodology	41
3.5.1	Dataset Preparation	41
3.5.1.1	Dataset D1:.....	41
3.5.1.2	Dataset D2:.....	42
3.5.1.3	Dataset D3:.....	44
3.5.2	Pre-processing for Proposed method:.....	45
3.5.2.1	Unequal Datasets:.....	45
3.5.3	Features Selection	46
3.5.4	Algorithms Evaluation.....	46
3.5.5	Result Compilation	47
3.5.5.1	Performance Metrics.....	47
3.5.5.2	Confusion Matrix.....	48
3.5.5.3	Accuracy.....	48
3.5.5.4	Precision	49
3.5.5.5	Recall	49
3.5.5.6	F-Measure	49

3.5.5.7	Cross-validation.....	50
3.6	Summary	50
Chapter 4	51
Result and Evaluation	51
4.1	Chapter Overview	51
4.2	Experiment Setup	51
4.3	Paper P1: Method 1	52
4.4	Paper P2: Method 2	54
4.5	Paper P3: Method 3	56
4.6	Paper P4: Method 4	58
4.7	Paper P5: Method 5	60
4.8	Paper P6: Method 6	62
4.9	Comparative Analysis (Best Optimal Method).....	64
4.10	Evaluation of NMPREDICTOR	67
4.11	Summary.....	71
Chapter 5	73
Conclusion and Future work	73
5.1	Chapter Overview	73
5.2	Thesis Summary	73
5.3	Future Work	75
REFERENCES	76

ABSTRACT

Vulnerabilities are known to be difficult to detect and prevent, especially in the context of web application. Although a significant research on web application security has been ongoing for a while, these applications have been a major source of problems and their security continues to be challenged. An important part of the problem derives from vulnerable source code of web applications. In order to overcome web vulnerabilities, different penetration tester used variety of techniques such as secure programming, static analysis, dynamic analysis, hybrid analysis and machine learning. Machine learning is consider an approach to prevent web vulnerabilities with a wide range of web applications because it is more preferable and does not have problems of false positive rate.

There are numerous method proposed for detecting web vulnerabilities based on machine learning. It is very difficult to measure, which method is efficient to secure web application. Furthermore, there is lack of study found that targets the comparison of machine-learning method to find out optimal method. However, comparative study is required to understand the path that could be followed by different penetration tester. In this thesis we use six different methods based on machine learning. In order to find optimal method for existing studies, decision were taken on Drupal metrics file with J48 and random forest. We have implemented NMPREDICTOR method with the feature extraction, performance parameters, classifiers with default parameters and 10k cross validation. Training data is passed through J48 and random forest to form a training model on which testing data is predicted and analyzed. Our results state that, to prevent web vulnerabilities VULPREDICTOR shows better results as compared to all others methods. We have found much higher accuracy of NMPREDICTOR method with respect to those reported by existing studies.