

C02047: Doctor of Philosophy
CRICOS Code: 058666A
33875 PhD Thesis: Computer Systems
March 2021

A Study on
Feature Analysis and Ensemble-based
Intrusion Detection Scheme using
CICIDS-2017 dataset

Upasana Nagar

School of Electrical and Data Engineering
Faculty of Engineering & IT
University of Technology Sydney
NSW - 2007, Australia

A Study on
Feature Analysis and Ensemble-based
Intrusion Detection Scheme using
CICIDS-2017 dataset

*A thesis submitted in partial fulfilment of the requirements
for the degree of*

Doctor of Philosophy
in
Computer Systems

by

Upasana Nagar

under the supervision of

Dr. Priyadarsi Nanda

to

School of Electrical and Data Engineering
Faculty of Engineering and Information Technology
University of Technology Sydney
NSW - 2007, Australia

March 2021

ABSTRACT

One of the primary security research challenges faced by traditional IDS methods is their inability to handle large volumes of network data and detect modern cyber-attacks with high detection accuracy and low false alarms. Hence, there is a need for efficient and reliable IDS schemes that can tackle this ever-changing cybersecurity paradigm. Machine learning techniques are hence, becoming very popular in designing modern intrusion detection systems. Several supervised and unsupervised machine learning techniques have been used in literature; however, the IDS classification efficiency is affected by noisy data in high dimensional datasets. The role of feature selection is significant as the feature selection process eliminates the redundant and noisy data and further selecting optimal feature subset enables reduction of high dimensional IDS datasets. Machine learning algorithms are extensively being used for intrusion detection. However, research has proved that the performance of multiple classifier-based IDS is far better than an IDS classifier, which has given us the motivation to develop an ensemble-based intrusion detection model. Lastly, the benchmark IDS datasets currently being used for the evaluation of IDS schemes are outdated and do not represent modern-day attacks. The CICIDS - 2017 dataset is offered by the University of New Brunswick. It is the latest publicly available dataset for intrusion detection. However, there are a significantly low number of research studies conducted using this dataset which also focus on optimal feature selection. This dataset has a good potential to be used as a future benchmark intrusion detection dataset as it covers the modern-day system setup and threat profile and the dependency on outdated IDS datasets can be removed. There is a need to benchmark the performance of modern IDS datasets using machine learning ensemble-based classifiers. This thesis aims to address the issues by proposing a new intrusion detection framework using ensemble-based feature selection method for generating a low dimensionality feature subset and ensemble-based intrusion detection framework to benchmark the performance of the CICIDS - 2017 dataset. The proposed scheme is beneficial for research community as it combines the use of the latest available IDS dataset with ensemble technique for feature selection and ensemble-based intrusion detection model.

AUTHOR'S DECLARATION

I, *Upasana Nagar* declare that this thesis, is submitted in fulfilment of the requirements for the award of *Doctor of Philosophy*, in the *School of Electrical and Data Engineering, FEIT* at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Production Note:
Signature removed prior to publication.
SIGNATURE: _____
[Upasana Nagar]

DATE: 30th March, 2021

DEDICATION

*I dedicate this achievement to the Supreme Almighty ...
For bestowing the blessing called "Life" ...*

*§
Immensely grateful for His guidance to "Turn Inwards" !!!

*I also dedicate this achievement to my parents ...
Papa & Mumma, you are my friend, philosopher and guide !!!!

ACKNOWLEDGMENTS

I am immensely grateful to the Almighty God for His kind "Grace". He has helped me conquer my fears with His guiding light, helping me achieve this Doctor of Philosophy Degree.

At this moment of accomplishment, I would like to express my deepest gratitude and appreciation to my supervisor Dr. Priyadarsi Nanda. This dissertation would not be possible without his able guidance and unwavering support. I am grateful for your expert feedback and valuable suggestions that helped me progress in my research. I am indebted to your kindness in offering your time and availability whenever I came to you for research discussion. Thank you for always being the extremely patient, positive, understanding, flexible and motivational guide, through out my research tenure. You gave me hope whenever I felt hopeless. Thank you for carving the researcher in me and showing me how to focus on seeking solutions to problems and always keeping a positive approach, while working on my research. I am blessed to have you as my mentor and my "Guru" for you have never given up on me even in times when I did give up on myself. I also wish to thank my co-supervisor Dr. Xiangjian He, for his valuable feedback and guidance during my research. Thank you for always offering kind and positive words of encouragement to keep me motivated in pursuing my research.

I would also like to thank my senior researchers Dr. Zhiyuan Tan, Dr. Aruna Jamdagni and Dr. Mohammed Ambu Saidi for their time and valuable guidance during the beginning of my research that helped me move ahead with my work. Thank you for all the knowledge sharing sessions that helped me build my research foundations. Many thanks to Eriyani and Aprillia from the School of Electrical and Data Engineering at the University of Technology Sydney for helping me with guidance with administrative tasks. I would also like to acknowledge Dr. Chandranath Adak for providing the PhD thesis template without which writing this thesis would have been an uphill cumbersome task. I would like to thank John Hazelton for his meticulous and prompt proof reading for my dissertation.

I am grateful for my time spent at UTS doing my research work as this gave me an opportunity to foster some great friendships. Big heartfelt thanks to Ashish Nanda and Amber Umair for always being there to inspire and motivate during my research. You both have helped me see the silver lining behind the dark clouds when I needed

it most. Thank you to Annie Baskaran, Nazar Waheed, Chau Nguyen, Nisha Malik, Madhumita Takalkar and Sara Farahmandian for their wonderful camaraderie, friendly banter, advice and support during my research and even outside. A very special thank you to my lovely friends Rattandeeep Kaur, Sai Kiran Tadepalli (Sai Garu) and Ramya Bati for their nurturing friendships.

I take this opportunity to thank my father Dr. T.N.Nagar, and my mother Mrs. Shobha Nagar, for their unconditional love, trust, support, encouragement, guidance and motivation in my life, especially as I encountered challenges during my research journey. You have taught me the power of prayer and built in me the resilient, "never give up" attitude towards life. You inculcated values of dedication, perseverance, sincerity, integrity, truthfulness and loyalty. These golden values have made me accomplish my goals and become who I am today. You have always reinforced the belief in me that I can move mountains if I put my mind to it! Thank you from the bottom of my heart for being my safety net and giving me the mental strength and confidence to realize my dreams into reality. I owe you my existence and I can not thank God enough for blessing me with the most awesome and caring parents. This PhD has been made possible only because of your presence and relentless support in the final year of my research. I would like to dedicate my thesis to you both.

Special thanks to my younger sister, Kankana Nagar, for being the playful, happy person and also for being the wiser one between us. You always brighten up my moods and lift my spirits with your presence.

I thank my mother-in-law, Hurshit Nagar, for her love, best wishes and blessings. I would also like to fondly remember my late father-in-law, Hursh Nagar, whose blessings are always with me. He would be so proud of my big achievement.

Last but certainly not least, I would also like to thank my incredible husband, Prasanavanam Nagar for believing in my dream and showering unconditional love and support during my research journey and in life. My heartfelt gratitude goes to my precious children and my forever cheerleaders, Aum and Anika. Their magical smiles always inspire me to keep giving my best in life.

Words fail to express my gratitude towards my extended family and all my family friends, who have always encouraged me to chase my research goals. I wish to conclude this acknowledgement, with a quote from one of my favourite books, "The Alchemist", by Paulo Coelho :

"And, when you want something, all the universe conspires in helping you to achieve it."

LIST OF PUBLICATIONS

1. Nagar, U, Nanda, P & He, X Feature Analysis and Ensemble-based Intrusion Detection Scheme using CICIDS - 2017 dataset. (Submitted to Wiley, Concurrency and Computation Practice and Experience Journal.)
2. Nanda, P, Arain, A & Nagar, U 2019, 'Network Packet Breach Detection using Cognitive Techniques', Smart Systems and IoT: Innovations in Computing, Second International Conference on Smart IoT Systems - Innovations in Computing (SSIC-2019), Springer, Jaipur, India.
3. Nagar, U, Nanda, P, He, X & Tan, Z 2017, 'A Framework for Data Security in Cloud using Collaborative Intrusion Detection Scheme', SIN'17: Proceedings of the 10th International Conference on Security of Information and Networks, International Conference On Security Of Information And Networks, ACM Digital Library, Jaipur, India, pp. 188-193.
4. Ambusaidi, MA, He, X, Tan, Z, Nanda, P, Lu, L & Nagar, U 2014, 'A novel feature selection approach for intrusion detection data classification', 2014 IEEE 13th International Conference on Trust, Security and Privacy in Computing and Communications, IEEE International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom), IEEE Computer Society, Beijing, pp. 82-89.
5. Tan, Z, Nagar, U, He, X, Nanda, P & Liu, R 2014, 'Enhancing Big Data Security with Collaborative Intrusion Detection', IEEE Cloud Computing Magazine, pp. 34-40.

TABLE OF CONTENTS

List of Publications	vi
List of Figures	x
List of Tables	xiii
1 Introduction	1
1.1 Background	1
1.1.1 Overview of Cyber threats and attacks	3
1.1.2 Intrusion Detection System - Need for robust cyber security	6
1.1.3 Research Motivation : Challenges with Intrusion Detection schemes, Feature selection techniques and IDS benchmark datasets	10
1.2 Problem Statement and Research Contribution	11
1.3 Research Overview : Objectives and Methodology	14
1.3.1 Research Objectives	14
1.3.2 Research Methodology	15
1.4 Thesis Organization	16
2 Intrusion Detection Systems - Taxonomy and Related Works	18
2.1 Taxonomy of Intrusion Detection Systems	18
2.1.1 Intrusion Detection System Based on Data sources	20
2.1.2 Intrusion Detection System Based on Detection Methodologies	23
2.1.3 Intrusion Detection System Based on System Structure	24
2.1.4 Intrusion Detection System Based on Audit Time of Data	25
2.1.5 Intrusion Detection System Based on Action taken after detection	26
2.1.6 IDS Based on Implementing Approaches	26
2.2 Taxonomy of Network Threats	29
2.3 Review of Intrusion Detection Datasets	35

2.4	Literature review of related works	39
2.4.1	Supervised Machine Learning Algorithms	43
2.4.2	Unsupervised Machine Learning Algorithms	48
2.4.3	Ensemble based Intrusion Detection system	48
2.5	Commonly used metrics employed for IDS evaluation	50
2.6	Chapter summary	54
3	Proposed Ensemble based Feature Selection approach	55
3.1	Further discussion on KDD 99, NSL-KDD and CICIDS-2017 datasets . .	56
3.1.1	KDD 99	58
3.1.2	NSL - KDD	59
3.1.3	UNSW - NB15	60
3.1.4	CICIDS - 2017	63
3.2	Feature Selection using Machine Learning techniques	68
3.3	Proposed Ensemble-based Feature Selection Approach using Machine learning methods	72
3.4	Chapter Summary	77
4	Proposed Ensemble based Machine Learning Technique for Intrusion Detection model	78
4.1	Ensemble-based Machine Learning approach for Network Anomaly Detec- tion - An Overview	79
4.1.1	Techniques to combine ensemble classifiers	80
4.2	Proposed Ensemble-based Network Anomaly Detection Intrusion Detec- tion Framework	86
4.2.1	Data Pre-processing	87
4.2.2	Proposed Ensemble-based Feature selection technique	89
4.2.3	Proposed ensemble-based IDS framework explained	89
4.3	Chapter Summary	93
5	Results and Analysis	95
5.1	Results from the implementation of Proposed Framework with CICIDS - 2017 Dataset	96
5.1.1	CICIDS - 2017 Data Filtering Process	96
5.1.2	CICIDS - 2017 Ensemble-based Feature Selection Process	97

5.1.3	Results from implementation of Proposed Ensemble-based intrusion detection model with CICIDS - 2017	100
5.2	Feature Selection and Ensemble-based Detection for NSL - KDD dataset .	103
5.3	Chapter Summary	106
6	Conclusion And Future works	109
6.1	Thesis Contributions	109
6.2	Future Work	111
A	Appendix	113
A.0.1	ROC graphs generated for Benign and attack classes of the CICIDS - 2017 dataset and the NSL - KDD dataset	113
	Bibliography	122

LIST OF FIGURES

FIGURE	Page
1.1 A generalized framework for a conventional cybersecurity system [57]	8
2.1 Taxonomy of Intrusion Detection Systems.	20
2.2 Comparison of IDS technology types based on their positioning within the computer system [90]	22
2.3 Centralized and distributed IDS (M=Monitors and A=Analysis units) [183] .	25
2.4 Taxonomy of network anomaly detection approaches [116]	28
2.5 Threat profile for cloud computing [85]	33
2.6 Distribution of network threats covered in [79]	34
2.7 Comparison of selected benchmark IDS Datasets [90]	39
2.8 Distribution of datasets for evaluating IDS performance in literature [79] . .	40
2.9 Cloud computing security issues [85]	41
2.10 Cloud computing security issues [173]	42
2.11 Collaborative Intrusion Detection Scheme for Public Cloud protect[122] . . .	43
2.12 Decision Tree representation [111]	44
2.13 Generic Algorithm flow of execution [111]	45
2.14 Various layers in Neural Network [111]	46
2.15 Linear SVM [111]	47
2.16 Distribution of Intrusion Detection Techniques based on research included in [79]	49
2.17 Confusion Matrix explained [156]	53
3.2 Data set record distribution KDD 99 and NSL-KDD [188]	58
3.1 Feature set of KDD 99 dataset [90]	59
3.3 Test bed set up for UNSW-NB15 [115]	60
3.4 UNSW-NB15 Basic Features [117]	61
3.5 Data set record distribution of UNSW-NB15 [117]	62

3.6	UNSW-NB15 Content Features [117]	62
3.7	UNSW-NB15 Time Features [117]	63
3.8	UNSW-NB15 Flow Features [117]	63
3.9	UNSW-NB15 Generic and Connection Features [117]	64
3.10	Test bed framework for CICIDS - 2017 dataset [153]	65
3.11	Data set record distribution of CICIDS - 2017 [163]	66
3.12	CICIDS 2017 Dataset Attack Category [188]	67
3.13	CICIDS 2017 Dataset Features [20]	68
3.14	Wrapper based Feature Selection Method [177]	69
3.15	Filter based Feature Selection Method [177]	70
3.16	Embedded Feature Selection Method[177]	71
3.17	Advantages and Disadvantages of Feature Selection Methods[177]	72
3.18	Proposed Ensemble-based Feature Selection approach	73
4.1	Distribution of Ensemble Base Classifier combining method [170]	81
4.2	Distribution of Combination Rule for Majority Voting using data from [170]	82
4.3	Trend of using Ensemble-based methods between 2015 - 2020 [170]	86
4.4	Proposed Ensemble Feature Selection and Ensemble-based Network Intrusion Detection Framework	88
4.5	Classification example of KNN for K=5 [90]	90
4.6	C4.5(J48) pseudo code [5]	91
5.1	Individual features generated for IG, CFS and PSO	99
5.2	Class distribution for NSL - KDD dataset	104
5.3	Proposed Feature Selection and Ensemble Model summary	107
A.1	CICIDS - 2017 Benign	113
A.2	CICIDS - 2017 Bot	114
A.3	CICIDS - 2017 BruteForce	114
A.4	CICIDS - 2017 DDoS	114
A.5	CICIDS - 2017 Hulk	115
A.6	CICIDS - 2017 FTP	115
A.7	CICIDS - 2017 SSH	115
A.8	CICIDS - 2017 GoldenEye	116
A.9	CICIDS - 2017 HeartBleed	116
A.10	CICIDS - 2017 Ilfiltration	116

A.11 CICIDS - 2017 Port Scan	117
A.12 CICIDS - 2017 SlowHTTP	117
A.13 CICIDS - 2017 Slowloris	117
A.14 CICIDS - 2017 SQL Injection	118
A.15 CICIDS - 2017 XSS	118
A.16 NSL - KDD U2R	119
A.17 NSL - KDD Normal	119
A.18 NSL - KDD Probe	120
A.19 NSL - KDD DDoS	120
A.20 NSL - KDD R2L	121

LIST OF TABLES

TABLE	Page
3.1 Comprehensive Summary of Publicly available IDS Datasets [176]	57
4.1 Comparison of Ensemble-based Network Anomaly IDS schemes mentioned in this work	84
5.1 Record details for CICIDS - 2017 after Data Filtering	97
5.2 CICIDS - 2017 Selected Ensemble Features after Data Filtering	98
5.3 Comparative results of Feature Selection methods for CICIDS - 2017	99
5.4 Overall Performance of Ensemble model with CICIDS2017 Dataset (68 features)	101
5.5 Overall Performance of Ensemble model with CICIDS2017 (17 features) . . .	101
5.6 Attack Class Performance of Ensemble model for CICIDS2017 (17 features) .	102
5.7 Comparison of proposed model with existing schemes using CICIDS - 2017 .	103
5.8 NSL - KDD Selected Ensemble Features	104
5.9 Overall Performance of Ensemble model with NSL-KDD dataset (41 features)	105
5.10 Overall Performance of Ensemble model with NSL-KDD dataset (15 features)	105
5.11 Attack Class Performance of Ensemble model with NSL-KDD dataset (15 features)	105
5.12 Overall Performance of Ensemble model with UNSW-NB15 dataset (44 features)	106
5.13 Overall Performance of Ensemble model with UNSW-NB15 dataset (24 features)	106
5.14 Comparison of the proposed scheme based on model generation time	106