# DETECTION AND CLASSIFICATION OF MALICIOUS NETWORK STREAMS IN HONEYNETS

A THESIS PRESENTED IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY
IN
COMPUTER SCIENCE
AT MASSEY UNIVERSITY, PALMERSTON NORTH,
NEW ZEALAND.

Fahim U.H. Abbasi

2013

# Abstract

Variants of malware and exploits are emerging on the global canvas at an ever-increasing rate. There is a need to automate their detection by observing their malicious footprints over network streams. Misuse-based intrusion detection systems alone cannot cope with the dynamic nature of the security threats faced today by organizations globally, nor can anomaly-based systems and models that rely solely on packet header information, without considering the payload or content.

In this thesis we approach intrusion detection as a classification problem and describe a system using exemplar-based learning to correctly classify known classes of malware and their variants, using supervised learning techniques, and detect novel or unknown classes using unsupervised learning techniques. This is facilitated by an exemplar selection algorithm that selects most suitable exemplars and their thresholds for any given class and a novelty detection algorithm and classification algorithm that is capable to detect, learn and classify unknown malicious streams into their respective novel classes. The similarity between malicious network streams is determined by a proposed technique that uses string and information-theoretic metrics to evaluate the relative similarity or level of maliciousness between different categories of malicious network streams. This is measured by quantifying sections of analogous information or entropy between incoming network streams and reference malicious samples. Honeynets are deployed to capture these malicious streams and create labelled datasets. Clustering and classification methods are used to cluster similar groups of streams from the datasets. This technique is then evaluated using a large dataset and the correctness of the classifier is verified by using "area under the receiver operating characteristic curves" (ROC AUC) measures across various string metric-based classifiers. Different clustering algorithms are also compared and evaluated on a large dataset.

The outcomes of this research can be applied to aid existing intrusion detection systems (IDS) to detect and classify known and unknown malicious network streams by utilizing information-theoretic and machine learning based approaches.

# Acknowledgements

All praises and thanks to Allah the exalted, the creator and lord of the universe, the source of all knowledge and wisdom. I thank Allah for all His blessings and especially for blessing me with good health, audacity, potential and ability to complete this manuscript.

I wish to express my profound gratitude to my supervisors Prof. Richard Harris, Dr. Giovanni Moretti and Prof. Stephen Marsland for their time, wisdom, patience, keen interest, inspiring guidance and valuable suggestions in planning, conducting and writing this manuscript. Thank you for sharing your wealth of knowledge with me and thank you for teaching me scientific methods and research. I would like to thank all my teachers who taught me in school, college and university in Pakistan and New Zealand. You all played your part that helped me to achieve my goals.

I have no words to express my gratitude to my father Mansoor-ul-Huda Abbasi and my mother Dr. Arsala Mansoor for all they have done for me. Thank you for believing in me and always encouraging me to pursue a Phd. I wish my father was alive and could be with us today to cherish this achievement. Both my parents were awarded a prestigious scholarship for Postgraduate studies overseas (in 1985) and I am glad to be their flag bearer. I would like to thank my siblings, my sister Asra, brother-in-law Irfan, and my brother Ibrahim. Thank you for your wishes and support.

I would like to thank my wife Rooshan, who has always stood firm by my side through good and difficult times during my PhD. Your constant love and support played a major role in completion of my PhD.

I would like to thank my lovely daughters Fatima and Sara. You have provided the perfect distraction that allowed me to keep my sanity and provided a good motivation to finish my thesis.

I would like to thank my fellow Postgrads at Massey University, especially my office mates Graham Jenson and Jevon Wright for the valuable discussions and distractions.

# Contents

# List of Tables

# List of Figures

xviii