MASSEY UNIVERSITY

# MAXIMISING THE EFFECTIVENESS OF THREAT RESPONSES USING DATA MINING: A PIRACY CASE STUDY

## *Seung Jun Lee*

**Master of Information Sciences**

Information Technology

*To my*

*Beloved grandmother and uncle who left me during the writing of this thesis,*

*I dedicate this thesis to them and may their souls rest in peace…*

*Dedication also goes to my beloved parents and my wife who have always been*

*very supportive with much love and providing their sincere support to me*

*during the entire research of this thesis.*

*I also dedicate this thesis to my beloved grandfather who left us*

*a few years ago who always gave lots of support and love to my beloved family…*

MAXIMISING THE EFFECTIVENESS OF THREAT RESPONSES USING

DATA MINING: A PIRACY CASE STUDY

This thesis is presented in partial fulfillment of the requirements for the degree of

Master of Information Sciences in
Information Technology

School of Engineering and Advanced Technology
At Massey University Albany, Auckland, New Zealand

By

Seung Jun Lee, PG Diploma in Information Sciences and
B Information Sciences (Massey University)

August 2015

# Acknowledgements

I would like to thank my supervisor, Professor Paul Watters for his sincere support, clear guidance, encouragement and strong belief in me during the entire process and research of this thesis. It was much privileged to have an invaluable opportunity working on this thesis project with his continuous support and helpful insight throughout the past year.

I am also very grateful to my dearest family, especially my parents who have always been very supportive to me with much love, endless support and patience. I also would like to thank my wife, Sily for providing me with much support, dedication and love to our family. My gratitude also goes to all the staff and faculty members of the School of Engineering and Advanced Technology at Massey University, Albany for their positive support and assistance over the past few years. I also would like to thank all my friends at Massey University who were shared lots of invaluable moments and great ideas that helped me to enhance and build more motivation during the past few years of studies.

# Abstract

Companies with limited budgets must decide how best to defend against threats. This thesis presents and develops a robust approach to grouping together threats which present the highest (and lowest) risk, using film piracy as a case study. Techniques like cluster analysis can be used effectively to group together sites based on a wide range of attributes, such as income earned per day and estimated worth. The attributes of high earning and low earning websites could also give some useful insight into policy options which might be effective in reducing earnings by pirate websites. For instance, are all low value sites based in a country with effective internet controls? One of the practical data mining techniques such as a decision tree or classification tree could help rightsholders to interpret these attributes.

The purpose of analysing the data in this thesis was to answer three main research questions in this thesis. It was found that, as predicted, there were two natural clusters of the most complained about sites (high income and low income). This means that rightsholders should focus their efforts and resources on only high income sites, and ignore the others.

It was also found that the main significant factors or key critical variables for separating high-income vs low-income rogue websites included daily page-views, number of internal and external links, social media shares (i.e. social network engagement) and element of the page structure, including HTML page and JavaScript sizes. Further research should investigate why these factors were important in driving website revenue higher. For instance, why is high revenue associated with smaller HTML pages and less JavaScript? Is it because the pages are simply faster to load? A similar pattern is observed with the number of links. These results could form a study looking into what attributes make e-commerce successful more broadly.

It is important to note that this was a preliminary study only looking at the Top 20 rogue websites basically suggested by Google Transparency Report (2015). Whilst these account for the majority of complaints, a different picture may emerge if we analysed more sites, and/or selected them based on different sets of criteria, such the time period, geographic location, content category (software versus movies, for example), and so on. Future research should also extend the clustering technique to other security domains.

# Table of Contents

# List of Figures

# List of Tables

# List of Abbreviations

ARFF          Attribute-Relation File Format

ASCII         American Standard Code for Information Interchange

BT            Behavioral Targeting

CDA          Communications Decency Act

COPA        Children Online Protection Act

COPPA      Children Online Privacy Protection Act

CPC          Cost-Per-Click

CPM         Cost-Per-Thousand-Impressions

CSS          Cascading Style Sheet

CSV          Comma Separated Values

CTR          Click-Through-Rate

DMCA       Digital Millennium Copyright Act

DMOZ      Directory Mozilla

EU            European Union

FTC          Federal Trade Commission

GUI          Graphical User Interface

HIC           High Income Cluster

HTML        Hypertext Markup Language

IP             Internet Protocol

IQR          Interquartile Range

IT             Information Technology

LIC           Low Income Cluster

LVF          Lower Visual Field

NA            Not-Applicable

OVA          Online Video Advertising

| | |
|---|---|
| P/E | Price-To-Earnings ratio |
| PPA | Pay-Per-Auction |
| SAS | Statistical Analysis System |
| SEO | Search Engine Optimisation |
| TRA | Theory of Reasoned Action |
| URL | Uniform Resource Locator |
| US | United States |
| USC | University of Southern California |
| UVF | Upper Visual Field |
| WEKA | Waikato Environment for Knowledge Analysis |
| WOT | Web of Trust |